

# Identifikacija spletno specifičnih kolokacij pogostega besedišča

Senja Pollak

Inštitut Jožef Stefan

*Slovenščina na spletu in v novih medijih, Ljubljana, 26. 11. 2015*

- Kolokacije: tipične sopojavitve besed
- Pari besed vs. baza in kolokator
- Frekvenčni vs. frazeološki pristopi
- Luščenje in preučevanje kolokacij za različne aplikacije:  
(leksikografija, terminografija, strojno prevajanje, generiranje naravnega jezika, klasifikacija dokumentov, poučevanje tujega jezika, luščenje informacij, analiza sentimenta)

- Korpus JANES: tviti, forumi, blogi, komentarji
- Preučevanje jezika spletnih uporabniških vsebin
  - ortografija (jz, cist, nevem, kwaaa)
  - skladnja (sva šle, bom bil)
  - leksikalna raven (neologizmi, pomenski premiki, aktualne tematike)
    - *všeček, virantovanje*
    - *sledilec, hud (dober), stisniti (datoteko)*
    - *rabiti (potrebovati)*
    - preučevanje na ravni enobesednih in večbesednih izrazov (kolokacije)

Leme:

- splošno slov. besedišče (dan, otrok)
- spletno pogostejše besedišče (tvit, blog)
- spletno specifično besedišče (tajmlajn, tviteraš)

Leme:

- splošno slov. besedišče (dan, otrok)  
(splošne leme, nove kolokacije, njihovi kolokatorji, pomenski premiki)
- spletno pogostejše besedišče (tvit, blog)
- spletno specifično besedišče (tajmlajn, tviteraš)

*avto, beseda, cerkev, človek, dogodek, energija, fant,  
glava, hiša, igralec, jezik, knjiga, ljubezen, mati, namen,  
oče, pesem, roka, sistem, šola, telo, ura, vas, zgodba, želja*

- samo luščenje kolokatorjev splošnega besedišča
- izdelan API za hitri izvoz kolokacijskih seznamov
- 150 lem
- specifično = ekskluzivno
- okno -1 (+1)
- ...

- korpusa  $S$  in  $R$  (v istem jeziku)
- izbor lem  $L$
- za vsako lemo  $l$  iz  $L$ : najdi  $K_l$  (  $\langle k(\text{olokator}), l(\text{ema}) \rangle$  ), specifični za en korpus
- cilj: seta kolokacij  $K_{LS}$  (in  $K_{LR}$ )  $\langle k, l \rangle$  specifična za  $S$  in  $R$

- korpusa *Janes* in *Kres* (v istem jeziku)
- izbor lem  $L$
- za vsako lemo  $l$  iz  $L$ : najdi  $K_l$  (  $\langle k(\text{olokator}), l(\text{ema}) \rangle$  ),  
specifični za en korpus
- cilj: seta kolokacij  $K_{LS}$  (in  $K_{LR}$ )  $\langle k, l \rangle$  specifična za  $S$  in  $R$



- korpusa *Janes* in *Kres* (v slovenščini)
- izbor lem  $L$
- za vsako lemo  $l$  iz  $L$ : najdi  $K_l$  (  $\langle k(\text{olokator}), l(\text{ema}) \rangle$  ),  
specifični za en korpus
- cilj: seta kolokacij  $K_{LS}$  (in  $K_{LR}$ )  $\langle k, l \rangle$  specifična za  $S$  in  $R$

- korpusa *Janes* in *Kres* (v slovenščini)
- izbor lem  $L$  (*150 "splošnih", najpogostejših lem*)
- za vsako lemo  $l$  iz  $L$ : najdi  $K_l$  ( $\langle k(\text{olokator}), l(\text{ema}) \rangle$ ), specifični za en korpus
- cilj: seta kolokacij  $K_{LS}$  (in  $K_{LR}$ )  $\langle k, l \rangle$  specifična za  $S$  in  $R$

- korpusa *Janes* in *Kres* (v slovenščini)
- izbor lem  $L$  (150 "splošnih", najpogostejših lem)
- za vsako lemo  $l$  iz  $L$ : najdi  $K_l$  (  $\langle k(\text{olokator}), l(\text{ema}) \rangle$  ),  
specifični(=ekskluz.) za *Janes*
- cilj: set kolokacij  $K_{LS}$ , specifičen(=ekskluz.) za *Janes*

- funkcija *kolokacije* orodja SketchEngine

- funkcija *kolokacije* orodja SketchEngine: slovenščina
- 150 lem
- izvoz kolokacij Kres, Janes: API
- okno -1 (+1), logDice (Rychlý): leks. kolok.
- frekv. k: 5, n: 3 (n: 10)

		Freq	logDice
<a href="#">p</a>   <a href="#">N</a>	knjižen	67	8.100
<a href="#">p</a>   <a href="#">N</a>	tečaj	79	7.758
<a href="#">p</a>   <a href="#">N</a>	profesorica	38	7.490
<a href="#">p</a>   <a href="#">N</a>	učiteljica	36	7.134
<a href="#">p</a>   <a href="#">N</a>	raba	40	7.129
<a href="#">p</a>   <a href="#">N</a>	pravilen	58	6.956
<a href="#">p</a>   <a href="#">N</a>	polomljen	24	6.829
<a href="#">p</a>   <a href="#">N</a>	pogovoren	24	6.817
<a href="#">p</a>   <a href="#">N</a>	znanje	65	6.579
<a href="#">p</a>   <a href="#">N</a>	učenje	31	6.557
<a href="#">p</a>   <a href="#">N</a>	zboren	16	6.432
<a href="#">p</a>   <a href="#">N</a>	pouk	20	6.367
<a href="#">p</a>   <a href="#">N</a>	učiti	51	6.323
<a href="#">p</a>   <a href="#">N</a>	učitelj	27	6.206
<a href="#">p</a>   <a href="#">N</a>	profesor	23	6.023
<a href="#">p</a>   <a href="#">N</a>	obvladati	19	5.613
<a href="#">p</a>   <a href="#">N</a>	naučiti	29	5.442
<a href="#">p</a>   <a href="#">N</a>	translate	8	5.381
<a href="#">p</a>   <a href="#">N</a>	poučevanje	7	5.175
<a href="#">p</a>   <a href="#">N</a>	inštrukcija	7	5.158
<a href="#">p</a>   <a href="#">N</a>	v	2,397	4.920
<a href="#">p</a>   <a href="#">N</a>	brezhiben	6	4.866
<a href="#">p</a>   <a href="#">N</a>	študirati	9	4.806
<a href="#">p</a>   <a href="#">N</a>	uporaba	21	4.734
<a href="#">p</a>   <a href="#">N</a>	poznavanje	6	4.652
<a href="#">p</a>   <a href="#">N</a>	lektorat	4	4.446
<a href="#">p</a>   <a href="#">N</a>	spakedran	4	4.443

Kolokacije (lema: <i>družina</i> )	
Kres	Janes
član družine	enostarševska družina
ribiška družina	član družine
kraljeva družina	primarna družina
lovska družina	kraljeva družina
ustvariti družino	mlada družina
rejniška družina	ogrožena družina
mlada družina	cela družina
kmečka družina	ustvariti družino
njegova družina	( <i>n</i> ) članska družina
svoja družina	romska družina

Kolokacije (lema: <i>družina</i> )	
Kres	Janes
član družine	enostarševska družina
ribiška družina	član družine
kraljeva družina	primarna družina
lovska družina	kraljeva družina
ustvariti družino	mlada družina
rejniška družina	ogrožena družina
mlada družina	cela družina
kmečka družina	ustvariti družino
njegova družina	( <i>n</i> ) članska družina
svoja družina	romska družina

Kolokacije (lema: <i>družina</i> )	
Kres	Janes
član družine	enostarševska družina
ribiška družina	član družine
kraljeva družina	primarna družina
lovska družina	kraljeva družina
ustvariti družino	mlada družina
rejniška družina	ogrožena družina
mlada družina	cela družina
kmečka družina	ustvariti družino
njegova družina	( <i>n</i> ) članska družina
svoja družina	romska družina

Janes:8  
199, 1,23/M

Kres:5,7  
63; 0,52/M



- avtomatska primerjava seznamov
- ekskluzivnost (samo na seznamu Janes, i.e.  $Kres < 3$ )
- izvoz kolokacij Kres, Janes: API

#### Nastavitve:

- min logDice
- min frek. 10
- < razpr. (“idf”): odst. lem, ki imajo kolok. na seznamu oz. št. kolokac. seznamov s kolok./kolok. seznamami vseh lem

- vrednost *razp.* kot filter
- splošne funkcijske besede: (*razp* > 0, 9: *a*, *ampak*, *brez*, *dober*, *en*, *glede*, *isti*, *kak*, *kako*, *kakšen*, *moj*, *morati*, *nek*, *nekaj*, *niti*, *njihov*, *noben*, *oz.*, *reči*, *svoj*, *tak*, *tisti*, *torej*, *tvoj*, *vaš*, *veliko*, *vsak*, *zaradi.*)
- nove rabe jezika, na ravni kolokatorja:
  - pogoste besede z opuščeni strešicami: *vec*
  - nestand. okrajšave: *slo*,
  - pogosti tematski izrazi (*političen*, *janšev*)
  - nest. raba (*rabiti*: *rabiti prostor*, *hrano*)
  - pom. premiki kolokatorjev (*hud*: *huda slika*, *hud igralec*)
  - razlike med lematizaroji dveh korpusov (*edini* vs. *edin*)., itd.

- vrednost *razp.* kot filter
- splošne funkcijske besede: (*razp* > 0, 9: *a*, *ampak*, *brez*, *dober*, *en*, *glede*, *isti*, *kak*, *kako*, *kakšen*, *moj*, *morati*, *nek*, *nekaj*, *niti*, *njihov*, *noben*, *oz.*, *reči*, *svoj*, *tak*, *tisti*, *torej*, *tvoj*, *vaš*, *veliko*, *vsak*, *zaradi*.)
- nove rabe jezika, na ravni kolokatorja:
  - pogoste besede z opuščeni strešicami: *vec*
  - nestand. okrajšave: *slo*,
  - pogosti tematski izrazi (*političen*, *janšev*)
  - nest. raba (*rabiti*: *rabiti prostor*, *hrano*)
  - pom. premiki (*hud*: *huda slika*, *hud igralec*)
  - razlike med orodji za predproc. dveh korpusov (*edini* vs. *edin*).

ZANIMIVE BAZE (leme)

Nastavitve	Št. izluščenih kolokacijskih kandidatov
brez omejitev (logDice>0)	11.148
logDice>3	4.183
logDice>3, frek>=10	2.928
logDice>3, frek>=10, <i>razp</i> <0.1	1.661

Nastavitve	Št. izluščenih kolokacijskih kandidatov
brez omejitev (logDice>0)	11.148
logDice>3	4.183
logDice>3, frek>=10	2.928
logDice>3, frek>=10, <i>razp</i> <0.1	1.661

## Vsi pogoji:

- *homoseksualna družina*
- *narcisistična družina*
- *ožja družina*
- *razdreti družino*
- *Festival družin*
- (*<razpasti, družina>, razpadla družina (oba korpusa)*)

Nastavitve	Št. izluščenih kolokacijskih kandidatov
brez omejitev (logDice>0)	11.148
logDice>3	4.183
logDice>3, frek>=10	2.928
logDice>3, frek>=10, <i>razp</i> <0.1	1.661

## Vsi pogoji:

- *homoseksualna družina*
- *narcisistična družina*
- *ožja družina*
- *razdreti družino*
- *Festival družin*
- (*<razpasti, družina>*, *razpadla družina* (oba korpusa))

## Brez filtra *razp.*:

- *Janševa družina*
- *Praznična Družina*
- *<celje, družina>* (*cela družina*)



Nastavitve	Št. izluščenih kolokacijskih kandidatov
brez omejitev (logDice>0)	11.148
logDice>3	4.183
logDice>3, frek>=10	2.928
logDice>3, frek>=10, <i>razp</i> <0.1	1.661

Le  $\log\text{Dice} > 3$  (ni omejitve frek.):

- vse prej naštete +
- *raznospolna družina, partnerjeva družina*
- *črkovna družina (grafika), modelna družina (avtomobilizem), imenska družina*
- *razdirati družino*
- *časopis Družina, priloga Družina, Prehrana družine*
- *grožnja družini*

Nastavitve	Št. izluščenih kolokacijskih kandidatov
brez omejitev (logDice>0)	11.148
logDice>3	4.183
logDice>3, frek>=10	2.928
logDice>3, frek>=10, <i>razp</i> <0.1	1.661

## Brez omejitev:

- vse prej naštete +
- *homoseksualna, oligarhična, čefurska, gejevska*
- zelo splošni kolokatorji *super* in *ok*
- funkc. besede, ki z lemo ne tvorijo sintaktičnih in semantičnih kolokacij:  
*potem, kdo, pač* in *sploh*

Lema	Kolokatorji
država	vitek, zadolžiti, rušiti, pravično, ugrabiti, koruptiven, zavoziti, gnili, pokrasti, ugrabljen, zajebati
jezik	muca, eksotičen, tuja, ang, venetski, šparati
mama	foto, yo, platišče, mreža, odkar, doječ, vreme, kurac, feltna
moški	hetera, napasti, feminizacija, kastracija
zgodovina	servisen, preverljiv, smetisce, spisati, retuširanje, brisanje, odpasti

- neformalni kolokatorji *nabaviti avto, nategovati ljudi, frej dan*
  
- aktualne tematike, novi pojmi
  - *zajebati, zavoziti, pokrasti državo; feminizacija moških, kastracija moških*
  - *nepremičninski zakon, privatizacija vode, varčevalna politika*
  - *transspolna oseba*
  - lastna imena: *Depala vas; Čezvesoljska zombi cerkev blaženega zvonjenja;*
  - tujejezične prvina: *gratis knjiga, rimejk filma*
  - *Facebook skupina, startup podjetje, wifi točka*

- novi pomeni: *brisanje zgodovine*

Telefon na tiho, SMS-ji, *brisanje zgodovine* na računalniku: jasni indici, da te vara ...  
*Brisanje zgodovine* na Googlu je zadnji teden ali dva postal svetovni hit.  
... and FB doesn't make it easy. A pozna kdo kšno hitrejšo rešitev za *brisanje zgodovine* ?
- terminologija: *evklidski prostor, rizičen odnos*
- idiomatika
  - *muca jezik papala (oz. popapala, popapcala), muca jezik papne*
  - *rana ura\*... , pitaj boga (pitati\*), jebo vas*

- izpuščanje strešic:
  - *smetisce zgodovine, vcerajsnja tekma, končna cena, soncen dan*
- napačna lematizacija:
  - *hetera moški, tuja jezik*
  - *anything glas*
  - *<foto, mama> Big foot mama*
  - *<kurac, mama> poln kurac mam oz. koji kurac mamo*
- nerazpršenost po korpusu:
  - *Pediater je odvetnik otroka je (300 pojavitev)*



- izpuščanje strešic: **rediakritizacija**
  - *smetisce zgodovine, vcerajsna tekma, koncna cena, soncen dan*
- napačna lematizacija: **izboljšanje lematizatorjev, problematična normalizacija**
  - *hetera moški, tuja jezik*
  - *anything glas*
  - *<foto, mama Big foot mama*
  - *<kurac, mama> poln kurac mam oz. koji kurac mamo*
- nerazpršenost po korpusu: **deduplikacija, račun razp. med uporabniki (korak po luščenju )**
  - *Pediater je odvetnik otroka je (300 pojavitev)*

- splošne leme
  
- specifične leme:
  - drag/ljub/znani/slovenski/pravi **tviteraš**
  - nova/huda/lepa/dobra **profilka**
  - meja/število **všečkov**
  - prava/velika **bizarka** nova/huda/lepa/dobra **profilka**

## Raziskava:

- luščenje kolokacij “splošnega” besedišča
- aktualne tem., novi pojmi, pomenski premiki
- tujejez. prvine, lastna imena
- napake luščenja

## Nadalnje delo:

- primerjava besednih skic (API!), kolokacije spletnega besedišča
- relevantnost kandidatov
  - za vključevanje v slovarje: slovar spletne slovenščine, kolokacijki slovar, SSSJ (Pollak, Arhar Holdt, Gantar)
  - kategorizacija (Pollak, Arhar Holdt, Gantar)
  - aktualne tematike, analiza diskurza (Škrjanec, Sobočan, Pollak)
  - frazeologija