

# **Luščenje in jezikoslovna analiza kolokacij iz korpusa Šolar**

TADEJA ROZMAN

ŠPELA ARHAR HOLDT

SENJA POLLAK

IZTOK KOSEM

JTDH 2016

# Okvir

## korpusna leksikalna analiza – usvajanje besedišča

- Špela Arhar Holdt in Tadeja Rozman. 2015. Možnosti uporabe podatkov iz korpusa Šolar za pripravo slovarskih priročnikov.
- Mojca Stritar Kučuk: Napake besedišča. (V: Izток Kosem idr. 2012.)

# Raziskovalno vprašanje

Kaj nam lahko korpus Šolar pove o usvajanju kolokacij  
in

kako te informacije uporabiti pri načrtovanju pouka leksike ter  
pripravi didaktičnih gradiv in jezikovnih priročnikov,  
namenjenih šolski populaciji?

# Raziskovalni problem

- specifičnosti v rabi kolokacij pri mladini v primerjavi z običajno rabo odraslih v splošni rabi
  - primerjalna analiza korpusov Šolar in Kres
- detekcija težav v rabi kolokacij pri mladini
  - analiza učiteljskih popravkov v Šolarju

# Primerjalno luščenje kolokacij: Šolar vs. Kres

- **Primerjalno luščenje kolokacij** ([Pollak in Arhar Holdt, 2015](#)):
  - za seznam besed v korpusu: avtomatski izvoz kolokacij iz orodja SketchEngine ([Kilgariff et al., 2004](#), [API Pollak, 2015](#)): frekvence, kolok. vrednosti (LogDice), povezava na zglede
  - izračun primerjalnih kolokacijskih vrednosti LD-diff (LDŠolar-LDKres), dodatni pogoji (bes. vrsta, minimalne vrednosti, ...)
- **Motivacija:** nove kolokacije, enake oz. primerljive kolokacije, tipične kolokacije vs. redke zveze
- **Izbor za luščenje in analizo:**
  - top 100 samost. v korpusu Šolar (*človek, življenje, ljubezen, otrok, čas* itd.),
  - pozicija pred lemo (in bes. vrsta: pridevnik, glagol ali samostalnik)
  - $\logDice \geq 3$  in relativna frekvenca  $\geq 2$  promila (vsaj 2 pojavitvi v Šolarju oz. vsaj 200 v Kresu).

# Izluščeni podatki za vzorec [pridevnik + čustvo]

Kolokator	Iztočnica	lowUP	LD_diff	logDiceSolar	Frek-Sola	logDiceKres	Frek-Kres	linkSola	linkKres	Pojavljanje
negativen	čustvo	LOWER	-9,089	0,000	0	9,089	263	nolink	http://solar	SAMO_KRES
ljubezenski	čustvo	LOWER	-0,580	7,159	3	7,738	93	http://sola	http://solar	V_OBEH
globok	čustvo	LOWER	0,010	7,685	4	7,677	108	http://sola	http://solar	V_OBEH
romantičen	čustvo	LOWER	0,670	6,712	2	6,040	26	http://sola	http://solar	V_OBEH
pozitiven	čustvo	LOWER	1,500	8,447	8	6,951	70	http://sola	http://solar	V_OBEH
poln	čustvo	LOWER	2,040	6,424	2	4,381	19	http://sola	http://solar	V_OBEH
močen	čustvo	LOWER	2,410	9,764	28	7,354	177	http://sola	http://solar	V_OBEH
osnoven	čustvo	LOWER	2,810	6,438	2	3,625	13	http://sola	http://solar	V_OBEH
lep	čustvo	LOWER	3,210	6,508	4	3,303	18	http://sola	http://solar	V_OBEH
intimen	čustvo	LOWER	3,550	7,433	3	3,884	5	http://sola	http://solar	V_OBEH
pomemben	čustvo	LOWER	4,500	5,701	2	1,200	5	http://sola	http://solar	V_OBEH
ustrezen	čustvo	LOWER	5,540	6,805	2	1,264	2	http://sola	http://solar	V_OBEH
psihičen	čustvo	LOWER	7,359	7,359	3	0,000	0	http://sola	nolink	SAMO_SOLAR

# [pridevnik + čustvo] glede na LD\_Diff



Gim, 4. letnik, NG

Berk je nezmožen **psihičnih čustev** nesvobodo-miselen, neprilagodljiv intelektualec, ki se je kot partizan boril proti Nemcem.

Strok, 4. letnik, MB

Sama dajem mogoče malo večji poudarek **psihičnim čustvom**, ki znajo biti nenadomestljiva, čeprav se zavedam, da je telesna ljubezen prav tako pomembna.

[pridevnik + samostalnik] prekrivno:  
1810 zvez izmed 2266 izluščeni (79,9  
%)

Prvih 40 glede na frekvenco v korpusu Šolar

šolski spis (frek. 390), šolska naloga, današnji čas,  
nezakonska mati, družbene razmere, osnovna šola,  
dobro življenje, naslednji dan, glavni junak, srednja  
šola, materni jezik, današnji svet, glavni lik, prosti čas  
dober prijatelj, prava ljubezen, športni dan, glavna  
oseba, izgubljeni sin, pravi prijatelj, posmrtno  
življenje, svetovna vojna, celo življenje, krščanska  
vera, velik vpliv, spolni odnos, organska potreba,  
ljubljen oseb, človekova usoda, kulturni dom,  
nezakonski otrok, očetova smrt, velik problem, mlado  
dekle, vsakdanje življenje, današnja družba, povodni  
mož, dolg čas, Cankarjev roman, pisna naloga (frek. 36)

temeljno besedišče  
vzgojna vloga slovenščine



# [pridevnik + samostalnik] samo v Kresu: 279 zvez izmed 2266 izluščeni (12,3 %)

## Prvih 40 glede na logDice

pravna oseba (logDice 11,176), nadaljnje besedilo, fizična oseba, nadzorni svet, gospodarska družba, združene države, svetovalna služba, obveščevalna služba, diplomatska naloga, prečiščeno besedilo, delniška družba, zemljiška knjiga, socialno delo, državni svet, občinski svet, odgovorna oseba, posebna potreba, negativno čustvo, lastninska pravica, študentski dom, prednostna naloga, leva stran, škodljivi vpliv, varnostna služba, desna stran, izredna razmera, uradna oseba, tajna služba, državni svet, investicijska družba, delovni čas, terensko delo, prihodnje leto, posredniška družba, zvezna država, ločeno mnenje, brezposelna oseba, mednarodni odnos, ustavna pravica, letalska družba (logDice 8,357)

medpredmet. povezava  
neumetnostna besedila  
neliterarni eseji  
šolski slovar

# [pridevnik + samostalnik] samo v Šolarju: 177 zvez izmed 2266 izluščeni (7,8 %)

## Prvih 40 glede na logDice

Antigonina sestra (logDice 10,034), Črtomirjev vojak,  
Kreonov sin, Rebulov roman, Kreonov zakon, Bronjin  
mož, Antigonin brat, Lojzkin starš, Kreonova žena,  
Descartesova misel, kiparska naloga\*, zavedni del\*,  
Kovačičev roman, Polikarpov odnos, Jazonova žena,  
Hamletova težava, mrtvaški pot, Polinejkova sestra,  
zardele oči, Lojzkine oči, aplikativni cilj\*, Polikarpovo  
dejanje, Salomin odnos, Gregorjevo mesto,  
Sofoklejevo delo, novodobni razum, Bubijevo dekle,  
Bogomilina želja, Bronjina želja, Odisejevo mnenje,  
Johanaanova glava, Johannova glava, prerokova  
glava, edini način, psihično čustvo, edina oseba,  
Ožbejev oče, rumeni junak, edini prijatelj (logDice 7,241)

# Kvalitativna analiza popravkov v korpusu Šolar

- **Luščenje kolokacij:**

- vsaj en del zveze mora biti popravljen
- jezikovni popravek mora biti uvrščen v tip 'napaka besedišča'

- **Izbor za luščenje in analizo:**

- pridevnik + samostalnik (278)
- samostalnik + samostalnik (271)
- glagol + samostalnik (274)

# [pridevnik + samostalnik] statistika

- skupaj 278 primerov
- šum: 52 primerov (18,7 %)
- relevantnih za analizo = 226 primerov (81,3 %):
  - popravek besedne zveze: 47 primerov (20,8 %); popravek ene besede v zvezi: 179 (79,2 %)
  - esej/spis: 213 (94,2 %), pisni izdelek (učna ura): 13 (5,8 %)
  - samo slovenščina

# [pridevnik + samostalnik]

- praviloma nizke frekvence

red prišla <i>nova</i>	<i>punca</i>	<i>sošolka</i> po imenu !
<i>bivša</i>   <i>bivše</i>	<i>punca</i>	<i>dekle</i> . Na vsak na
e stvari. <i>Mlada</i>	<i>punca</i>	<i>dekleta</i>   <i>dekleta</i>
<i>mlada</i>   <i>mlado</i>	<i>punca</i>	<i>dekle</i> nadgrajeval
je o <i>najstniški</i>	<i>punci</i>	<i>dekletu</i> Fatimi Isn
<i>njo</i>   <i>prejšnje</i>	<i>punco</i>	<i>dekle</i> . Nisi ji hote
i   , kako <i>lepo</i>	<i>punco</i>	<i>dekle</i> ima. Julija t

<i>materin</i>	<i>materni jezik</i> in se
<i>materin</i>	<i>materni jezik</i> . Prer
<i>materin</i>	<i>materni jezik</i> lep in
<i>materin</i>	<i>materni jezik</i> . Razl
<i>materinega</i>	<i>maternega jezika</i> .
<i>materinega</i>	<i>maternega jezika</i> .
<i>materinskega</i>	<i>maternega jezika</i> .
<i>materinskem</i>	<i>maternem jeziku</i> in
<i>materinskem</i>	<i>maternem jeziku</i>
<i>materinskemu</i>	<i>maternemu jeziku</i>
<i>materinski</i>	<i>materni jezik</i> upora
<i>materinski</i>	<i>materni jezik</i> . Zave
<i>materinski</i>	<i>materni jezik</i> in to j
<i>materinski</i>	<i>materni jezik</i> . V ce
<i>materinski</i>	<i>materni jezik</i> in to j
<i>maternim</i>	<i>maternim jezik</i>   je

# [pridevnik + samostalnik] nerelevantni primeri

- vsebinsko narobe, besedilni popravki, nerazumljivi popravki, oblika

<p> Tartuffe je bil duhovnik, ki je odšel k **bogatemu** o se mama, ata in Tone odpravili na izlet v **Kranjsko** olj znano mislijo razsvetljenstva, ki jo je zapisal znan i upogibalo s tokom. </p> <p> Kreona, **Antigoninega** delje v Dobjo vas | , kjer imam pevske vaje. **Pevske mesto** meri 275 km<sup>2</sup>, kar nam pove, da je **največje im** . Ko **Polonij pove** | **je Polonij povedal** Klavdiju o a stoji \$\$\$ **v** | **na** nekem nam neznanem **prevoznem** največkrat opazimo športne copate. Je **velik** | **velika**

**Argonu** | **Oregonu** na dom in se naprav  
**goro** | **Rogaško Slatino** . Tam je vide  
**grški** | **francoski pisatelj** Descartes  
**očeta** | **strica** , | pa težko uvrstimo |  
**vaje** | **Te** trajajo do šestih , | in ko  
**mesto** | **središče** v Sloveniji. V Ljublj  
**Hamletovi** | **prinčevi ljubezni** do Ofelije,  
**sredstvu** | **avtobusu** , ki bi lahko bil avtc  
**ljubitelj** | **ljubiteljica** narave in živali. 9

# [pridevnik + samostalnik] skupine podobnih primerov leksikalnih popravkov

- tipologija popravkov
- tipologija besedišča
  - izrazi s področja slovenščine – ostalo besedišče (pomenska/pojmovna polja, abstraktnost, strokovnost, pogostost rabe ...)
- beseda – besedna zveza – mejni primeri

# [pridevnik + samostalnik] zaznamovano

- večinoma popravki besed, manj zvez

pedofilija, droge, nezvestoba pa so največkrat <b>rizični</b>	<b>faktorji</b>	<b>dejavniki</b> , da se starši od otrok umaknejo, stika z njimi ne
<b>jaz</b>   <b>Jaz</b> sem se čisto malo ,   – imam <b>nova</b>   <b>novo</b>	<b>frizuro</b>	<b>pričesko</b> in nekaj lažja sem. Čez dva tedna grem na morje s
o. Mislim, da ji je na začetku <b>pomembna</b>   <b>pomemben</b>	<b>funkcija</b>	<b>položaj</b> kraljice, ker ne ve, kaj se je v resnici zgodilo z
Ljudje smo bitja odnosov. Za naše <b>normalno</b>	<b>funkcioniranje</b>	<b>delovanje</b> potrebujemo drug drugega. Sami pripomoremo k c
e. Romeo je tik pred tem, da spije strup, a v <b>zadnjem</b>	<b>momentu</b>	<b>trenutku</b> do vežice pride duhovnik Lorenzo, ki Romea sezna
<b>rali</b> prislužiti kruh. Vzgajala jih je <b>kot</b>   <b>za pogumne</b>	<b>možake</b>	<b>ljudi</b> . To vidimo tudi v besedilu kot je: ”   „ Zdaj greš v bo
gremo skupaj na policijo. Ko smo prispeli na <b>najbližjo</b>	<b>policijo</b>	<b>policijsko postajo</b> , smo poiskali dežurnega policista in mu
. dolgočasna. Prav gotovo ga tudi to potre in je njegovo	<b>brezvezno</b>	<b>nesmiselno življenje</b> tako še bolj temno in pesimistično. V
e prikupiti s tem, da ga je zavajala in ga peljala na vsa	<b>živa</b>	<b>mogoča potovanja</b>   , misleč, da gre na službeno potovanje



# [pridevnik + samostalnik] nenapake

- bolj specifično (vsebinsko natančnejše, pomensko ožje)

<i>oddigrali</i>   <i>odigrali</i> predstavo o Romeu in Juliji, ki je bila <i>prava</i>	<i>drama</i>	<i>tragedija</i> .	<i>Meni kot osebi</i>
v različnem času, kar ni bilo značilno za <i>Antično</i>   <i>antično</i>	<i>gledališče</i>	<i>dramatiko</i> .	Tudi v življe
obdobju renesanse. <i>Romeo in Julija</i> sta <i>glavna</i>   <i>glavni</i>	<i>junaka</i>	<i>dram. osebi</i> .	Capuletovi so priredi
bo dobro godilo, če ga dobi s kakšno flašo. <i>Celotna</i>	<i>knjiga</i>	<i>drama</i> je dokaj grozna, saj se doga	
<i>so</i>   <i>sta</i> izpostavljeni dve najbolj temeljni sporočili <i>celotne</i>	<i>knjige</i>	<i>komedije</i> . Prvo sporočilo je povez	
jaz mislim. Odlomek Krst pri Savici je epsko- <i>lirska</i>	<i>pesem</i>	<i>pesnitev</i> , ker <i>sestavlja</i>   <i>vsebuje</i>	
tudi otroka ne. Delo Krst pri Savici je epsko- <i>lirska</i>	<i>pesem</i>	<i>pesnitev</i> , ki izpoveduje čustva in g	
<i>romantike</i> , ki traja od 1830 do 1848. Delo je epsko- <i>lirska</i>	<i>pesem</i>	<i>pesnitev</i>   , zgrajena iz treh delov	
oba potrudita za to. Nezakonska mati je <i>lirska</i>   <i>lirsko</i>	<i>pesnitev</i>	<i>delo</i> , ki spada v obdobje <i>romantik</i>	
<i>pisatelj</i>	<i>pisatelj</i>	<i>dramatik</i> in najboljši pisec tragedij	
<i>pisatelja</i>	<i>pisatelja</i>	<i>dramatika</i> Williama <i>Shakespeara</i>	
Temu sonetu <i>je sledil</i>   <i>sledi</i> Uvod, kjer govori o <i>sami</i>   <i>samih</i>	<i>vojni</i>	<i>bojih</i>   <i>bojih med</i>   <i>za</i> pogansko	
<i>dramatika</i> Williama Shakespeara je najbolj poznana <i>ljubezenska</i>	<i>zgodba</i>	<i>drama/</i> . Je dramsko delo, zgodba	
ipravljen storiti marsikaj. Njen značaj spoznavamo skozi <i>celotno</i>	<i>zgodbo</i>	<i>dramo</i> in če sprva deluje plehko in	

# [pridevnik + samostalnik] književnost in jezik

atik Sofokles in   <i>je</i> tragično	<b>dramatsko</b>	<i>dramsko besedilo</i> .
zvrst oz. odlomek je epsko in	<b>ilirsko</b>	<i>lirsko besedilo</i> , to
jena iz treh delov. Prvi del je	<b>posvetni</b>	<i>posvetilni sonet</i> Mat
.nes imamo vse preveč tujk in	<b>prenesenih</b>	<i>prevzetih besed</i> iz r
<i>večinoma</i> uporabljamo sleng,	<b>privzete</b>	<i>prevzete besede</i> , p
nska mati spada v liriko in je	<b>vložena</b>	<i>vložna pesem</i> . Nast
nska mati spada v liriko in je	<b>vložena</b>	<i>vložna pesem</i> . Nast

# [pridevnik + samostalnik] oblikovno blizu, pomensko različne

raje rekli država ponekod še vedno obnaša na tak **diskriminanten** | *diskriminatoren način* . Na  
Po končani osnovni šoli , | se je vpisal v srednjo **električno** | *elektro šolo* \$\$\$ – TOK, v \

e in junakovega ideala. </p><p> V delu je veliko **izpovednih** | *izpovedanih čustev* , še po  
*šestdeset* let. </p><p> Na našo šolo je prišel kot **književni** | *knjižni gost* , zato , | ker

velik nesmisel, karikira pomembne osebnosti. Za **Majhne** | *male ljudi* , kamor spada t  
” (umobolnici), prav tako tudi **veliki** | *starejši* in **mali** | *mlajši izvedenec* za metaf

Kreonu privedli Antigono, **kot** | storilko **tega** | **protizakonskega** | *protizakonitega dejanja* .  
ka je bila poražena. Da je bilo to delo Prešernova **osebna** | *osebna izpoved* , bi lahko  
Prešeren preusmeri slog predvsem na Črtomirovo **osebno** | *osebno raven* . Zato pravin  
cela družina ostane sama? Spomnim se družine iz **sosedovega** | *sosednjega kraja* , kjer se  
rja , | pridejo mladi in stari v Podpeco in gradijo **snežene** | *snežne gradove* . Imenujer  
<ljuček bi lahko neskončne ure pisali, kaj je prav, **trda** | *trdna načela* ali prilagajan

# [pridevnik + samostalnik] oblikovno blizu, pomensko različne

Gregor je partizan, ki na *rutinskem* **sprehodu** | *obhodu* opazi mladenko, ki s svojimi kre  
ila Črtomirja. Krščanska vera govori o **posvetnem** | *posmrtnem življenju* , | – da bosta lal  
ljubil Bogomilo in je hotel biti tudi v **posvetnem** | *posmrtnem življenju* znjo. Črtomir se j  
opazovalec | , saj je opazil čisto vse **možne** | *mogoče podrobnosti* . Od ženice do nek  
i. Poznam en primer ki zaradi **verske** **izpovednosti** | *verske pripadnosti/izpovedi* starši svoji l

# [pridevnik + samostalnik] kolokacije

- uporaba pomensko neustrezne kolokacije:

**posvetno** > posmrtno življenje, iz socialnega **reda** > sloja, na družbeni **ravni** > lestvici, na višji **ravni** > na višjem položaju, filmski **trak** > filmsko platno

- manj običajna kombinacija besed:

**nabrano** > naropano bogastvo, imeti **v popolni kontroli** > pod kontrolo, **zadelana** > zamašena ušesa, **zlati** > velik zaslužek, biti vzkipljive **jeze** > narave, vse **mogoče** > možne podrobnosti, **velika** > visoka samopodoba

- mešanje:

verska **izpovednost** > pripadnost/izpoved, nadmorska **gladina** > višina

# Ocena kvalitativne analize

- razmeroma zamudna
- frekvence relevantnih primerov so nizke, a so podatki pomembni > prispevajo k razumevanju procesov usvajanja „aktivnega“ besedišča > načrtovanje pouka leksike, izbor gradiv za usvajanje besedišča in snovanje priročnikov
- na podlagi trenutnih delnih analiz posploševanje rezultatov še ni mogoče oz. zaključki o procesih usvajanja kolokacij (leksike) ne morejo biti zanesljivi

# Kako naprej?

- primerjalno luščenje kolokacij razširiti tudi na manj frekventne samostalnice > primerjava s podatki iz kvalitativne analize
- primerjalno luščenje kolokacij uporabiti na verziji korpusa Šolar brez popravkov (šum!) > statistična analiza
- obe metodi (primerjalno luščenje kolokacij in kvalitativno analizo popravkov) razširiti na druge zveze
- ? analiza po razredih
- ? leksikalne analiza besed (kot osnova)
- ??? nove ideje ???

# Reference

Špela Arhar Holdt in Tadeja Rozman. 2015. Možnosti uporabe podatkov iz korpusa Šolar za pripravo slovarskih priročnikov. V: M. Smolej, ur., *Obdobja 34: Slovnica in slovar - aktualni jezikovni opis*, 1. del, str. 67–74. Znanstvena založba Filozofske fakultete UL. [http://centerslo.si/wp-content/uploads/2015/11/34\\_1-Arhar-Hol-Roz.pdf](http://centerslo.si/wp-content/uploads/2015/11/34_1-Arhar-Hol-Roz.pdf)

Vojko Gorjanc, Polona Gantar, Iztok Kosem in Simon Krek, ur. 2015. *Slovar sodobne slovenščine: problemi in rešitve*. Znanstvena založba Filozofske fakultete UL.

Adam Kilgarriff, Pavel Rychly, Pavel Smrz in David Tugwell. 2004. The Sketch Engine. V: G. Williams in S. Vessier, ur., *Proceedings of the Eleventh EURALEX International Congress, EURALEX 2004 Lorient, France July 6–10, 2004*, str. 105–116. Universite de Bretagne-sud.

Iztok Kosem, Mojca Stritar Kučuk, Sara Može, Ana Zwitter Vitez, Špela Arhar Holdt in Tadeja Rozman. 2012. *Analiza jezikovnih težav učencev: korpusni pristop*. Trojina, zavod za uporabno slovenistiko.

Nataša, Logar, Miha Grčar, Marko Brakus, Tomaž Erjavec, Špela Arhar Holdt, Simon Krek. 2012. *Korpusi slovenskega jezika Gigafida, KRES, ccGigafida in ccKRES: gradnja, vsebina, uporaba*. Trojina, Fakulteta za družbene vede.

Senja Pollak. 2015. Luščenje kolokacij iz korpusa uporabniških spletnih vsebin. V: M. Smolej, ur., *Obdobja 34: Slovnica in slovar - aktualni jezikovni opis*, 2. del, str. 601–607. Znanstvena založba Filozofske fakultete UL. [http://centerslo.si/wp-content/uploads/2015/11/34\\_2-Pollak.pdf](http://centerslo.si/wp-content/uploads/2015/11/34_2-Pollak.pdf)

Senja Pollak in Špela Arhar Holdt. 2015. Identifying corpus-specific collocations: the case of spoken Slovene. V: K. Gajdošová in A. Žáková, ur., *Natural language processing, corpus linguistics, lexicography: proceedings*, S. I., str. 117–125. RAM-Verlag.

Tadeja Rozman, Iztok Kosem, Nataša Pirih Svetina in Ina Ferbežar. 2015. Slovarji in učenje slovenščine. V: V. Gorjanc, P. Gantar, I. Kosem in S. Krek, ur., *Slovar sodobne slovenščine: problemi in rešitve*, str. 150–167. Znanstvena založba Filozofske fakultete UL.

Tadeja Rozman, Irena Krapš Vodopivec, Mojca Stritar Kučuk in Iztok Kosem. 2012. *Empirični pogled na pouk slovenskega jezika*. Trojina, zavod za uporabno slovenistiko.

Marko Stabej. 2011. Jezikovni potrošnik in potrošnica. *Sodobna pedagogika*, 62=128(2), 102–113. Zveza društev pedagoških delavcev Slovenije. <http://www.dlib.si/details/URN:NBN:SI:doc-RF8JNPLQ>