BESEDNE SKICE ZA SLOVENŠČINO. KRITIČNI POGLED

Olga Yeroshina Pobirk, Petra Zaranšek, Simon Šuster

Trojina, zavod za uporabno slovenistiko, Škofja Loka

UDK 811.163.6'374:811.163.6'322.3'324

Prispevek obravnava ozadje delovanja Besedne skice, ki je ena temeljnih funkcij orodja Sketch Engine, z vidika leksikografskega dela. Avtorji poleg opisa besedne skice kritično presojajo vlogo, ki jo imajo pri izdelavi in prikazu besednih skic označenost korpusa, njegova besedilna sestava, slovnična razmerja, določena v slovnici besednih skic, in nenazadnje tudi nastavitve pri prikazu kolokatorjev.

besedna skica, označenost korpusa, sestava korpusa, slovnična razmerja, statistični izračuni

The article discusses the functioning of Word Sketch, one of the basic functions of the Sketch Engine webbased application, particularly with regard to lexicographic work. In addition to describing an individual example, the authors assess the role that other factors play in the creation and display of Word Sketches, such as corpus annotation, corpus composition, the grammatical relations that are determined in the sketch grammar, and the settings regarding the arrangement of the displayed collocates.

Word Sketch, corpus annotation, corpus composition, grammatical relations, statistical scores

1 Uvod

Uporaba korpusov za opis jezika v leksikografske in druge raziskovalne namene tudi v slovenskem prostoru ni več novost, kar potrjuje vedno več študij, ki obravnavajo to tematiko. Korpus nudi obilico jezikovnih podatkov, vendar potrebujemo za njihovo obdelavo in interpretacijo zmogljiva ter »pametna« orodja. Primer takšnega sistema je večfunkcionalno orodje Sketch Engine (v nadaljevanju: SkE), ki omogoča različne prikaze in načine obdelave podatkov v jezikovnem korpusu.

SkE je orodje za prikaz in analizo jezikovnih podatkov v besedilnem korpusu, ki obsega naslednje osnovne funkcije, predvidene za vsebinsko delo z izbranim korpusom: Konkordance (*Concordance*), Seznam besed (*Word List*), Besedna skica (*Word Sketch*), Tezaver (*Thesaurus*), Primerjalna skica (*Sketch-Diff*). Od naštetih sodita med dve bistveni funkciji Konkordance in Besedna skica, ki je tudi predmet tega prispevka. V nadaljevanju bo najprej na kratko predstavljeno delovanje in namen besednih skic, nato pa bomo osvetlili ozadje njihovega delovanja, ki ga je treba vsaj do določene mere poznati, da znamo pravilno interpretirati dobljene podatke.

2 Besedna skica

Besedna skica predstavlja velik korak v smislu pomoči pri leksikografskem delu, odkar so prišli v rabo korpusi. Besedna skica prikazuje leksikalni profil izbrane iztočnice s podatki o njenem tipičnem sobesedilnem okolju (Gantar idr. 2009: 33). Gre za avtomatske, na korpusu temelječe sežetke slovničnega in kolokacijskega vedenja neke besede (Krek, Kilgarriff 2006), ki leksikografu dajejo hiter pregled nad dogajanjem okoli iztočnice. Leksikografu na ta način pomagajo tako pri pomenskem

ogeni

razdvoumljanju kot pri odločanju o tem, katere informacije naj tudi neposredno predstavi v slovarju oz. bazi, ki ga/jo sestavlja (Atkins, Rundell 2008: 110).

2.1 Sestava besedne skice

Slika 1 prikazuje del besedne skice za lemo *ogeni*.

Številka 39.083 nad besedno skico pove število pojavitev iskane leme, tj. besede v vseh njenih oblikah, v izbranem korpusu. Iz prvega stolpca z oznako *a_modifier* je mogoče razbrati, da je iztočnica *ogenj* v korpusu Fida-PLUS 8.961-krat rabljena kot samostalniško jedro v kombinaciji s pridevnikom. Spodaj sledi seznam pridevniških kolokatorjev, ki jo najpogosteje modificirajo in z njo tvorijo samostalniško besedno zvezo (npr. *taborni ogenj*). Stolpci *prec_z-d*, *prec_na-d* in *prec_v-d* prikazujejo besede, ki se običajno pojavljajo

Fida PLUS 620m freq = 39083

pred predlogi z, na, v, ki skupaj z iztočnico tvorijo predložno besedno zvezo (npr. odstaviti z ognja). V stolpcu gen_2 so našteti samostalniški kolokatorji, ki skupaj z iztočnico v rodilniku tvorijo samostalniško besedno zvezo (npr. prekinitev ognja). Stolpec z oznako is obj4 vsebuje glagole, ob katerih se ogenj pojavlja kot samostalnik v tožilniku (npr. zakuriti ogeni). Podobno stolpec is subj navaja glagole, ki so najpogosteje rabljeni ob iztočnici v imenovalniku, a ker je zgornja meja absolutne pogostosti v želji po bolj strnjeni besedni skici nastavljena na 50 pojavitev, imamo le en glagol, ki presega to mejo - goreti (ogenj gori). V stolpcu prec_prep so našteti najpogostejši predlogi, ki se pojavljajo ob iztočnici ogenj.

S klikom na »podčrtano številko« (Slika 1), je mogoče dostopati neposredno do konkordanc za izbrani kolokator. Številčne vrednosti desno od kolokatorjev pomenijo naslednje:

a modifier 8961 1.3 prec z-d 2208 5.0 prec v-d 2769 3.0 prec prep 10019 2.7 3074 1.8 580 96.24 Odstaviti 474 75.87 Zelezo 235 65.21 E z 2552 34.13 prekinitev taboren 706 81.98 navzkrižen <u>487</u> 83.4 53 34.38 noka 360 47.81 ob 509 30.09 Druhalec 125 81.19 igranje 138 32.5 zmeren 491 68.16 igra igra znajti 154 35.15 nad nad 212 29.71 požiralec 108 73.22 šibek 557 65.35 igrati 194 32.1 vreči 99 33.3 v 3115 28.12 prasketanje 61 60.19 blag 227 56.11 vzeti 128 30.11 poslati 157 30.26 okrog 70 25.98 prižiganje 71 56.05 Doj 64 29.59 pred pred 324 24.69 E kurjenje olimpijski 460 51.53 56 26.96 pošiljati 81 52.74 topniški <u>82</u> 50.9 umakniti 51 25.55 biti 267 19.4 proti 165 23.86 Druhanje <u>66</u> 49.15 93 10.34 169 11.02 na na 1547 23.7 ustavitev majhen 784 47.4 biti imeti 62 43.49 peklenski 102 47.04 >> okoli ol 53 19.98 🔲 vladavina 65 40.9 50 46.92 skozi 56 18.64 bog 71 37.23 stražen 300 44.87 prec na-d 2628 4.9 is obj4 3185 2.9 kot odprt 198 18.12 element 82 33.14 olje 776 70.79 zakuriti 279 71.92 brez močen 398 42.43 94 15.53 >> 335 54.83 zanetiti 320 70.05 večen 124 38.96 kuhati zaradi 89 13.46 80 50.6 411 69.47 kozica pogasiti 1557 1.8 počasen 76 33.58 iz 199 13.4 is subj 62 36.86 71 39.89 dušiti bruhati 154 56.44 od 🗎 goreti 72 32.93 144 13.34 prijateljski minuta 60 28.91 podtakniti 163 55.19 >> 75 27.42 k 77 13.06 svet 111 24.58 prižgati postaviti 218 52.01 54 24.98 0 vroč 73 5.49 59 5.6 kuriti 95 49.31 pri pri 55 24.0 <u>56</u> 3.53 Živ >> gasiti 53 39.8 za nizek nizek 57 19.34 179 2.89 ugasniti ugasniti 62 35.4 ро 71 2.37 notranji 56 18.95 160 34.02 zmanjšati srednji 51 18.18 >> opaziti 53 20.18 velik 192 15.75 >>

Slika 1: Del besedne skice za lemo ogenj.

»poděrtana številka« pove absolutno pogostost kolokatorja ob izbrani iztočnici, številka desno od nje pa izpostavljenost posameznega kolokatorja. Izpostavljenost (ang. salience score) je statistična vrednost, ki daje podatek o tem, kako močna je vez med kolokatorjem in izbrano iztočnico – večja ko je vrednost izpostavljenosti, močnejša je vez (prav tam).¹ Kolokatorji so v našem primeru razvrščeni glede na izračun izpostavljenosti in ne po absolutni pogostosti, zato se npr. sicer dva izmed najpogostejših glagolov biti in imeti v stolpcu prec_v-d nahajata na zadnjem mestu.

Pri razbiranju in interpretaciji podatkov, pridobljenih z besedno skico, je poleg poznavanja same uporabe orodja zelo pomembno tudi, da uporabnik pozna ozadje njegovega delovanja, zato bodo v nadaljevanju predstavljeni dejavniki, ki odločilno vplivajo na dobljene rezultate: besedilna sestava korpusa, naloženega v SkE, označenost korpusa in slovnična razmerja, definirana v datoteki s slovničnimi razmerji.

3 Ozadje delovanja besednih skic

3.1 Slovnica besednih skic

Za funkcije Besedna skica, Primerjalna skica in Tezaver potrebujemo slovnico besednih skic (ang. Sketch Grammar, dalje SBS), ki je nabor slovničnih razmerij, za katera želimo, da so prikazana v besednih skicah. Ta razmerja so definirana glede na zakonitosti v konkretnem jeziku, namen jezikovnih poizvedovanj ter nabor oznak, uporabljenih pri označevanju posameznega korpusa. Slovnico, ki je osnova za izdelavo besednih skic za slovenščino na 621-milijonskem korpusu FidaPLUS, je leta 2006 napisal Simon Krek (Krek, Kilgarriff 2006).

Kot primer slovničnega razmerja za korpus FidaPLUS navajamo recipročno zvezo pridevnika in samostalnika:

*DUAL =a modifier/modifies

2:[tag="P.*"] [tag="P.*"|word=","|word="se" |word="si"][0,5] 1:[tag="S.*"]

Razmerje je sestavljeno iz elementov, ki se začnejo in končajo z oglatimi oklepaji [], predstavljajo pa potencialne pojavnice v realnem besedilu. Znotraj elementov in med njimi so uporabljeni znaki iz regularnih izrazov :.," | *{} v kombinaciji z oblikoskladenjskimi oznakami, uvedenimi z atributom tag, ali dejanskimi koščki besedila, uvedenimi z atributom word. Zgoraj opisano razmerje predvideva naslednjo situacijo: na prvem mestu se mora pojaviti pridevnik v katerikoli obliki, sledi mu od nič do pet mest, na katerih se lahko pojavijo pridevnik, vejica, se ali si, na koncu pa mora biti samostalnik, ponovno v katerikoli obliki. S številkama 1 in 2 označimo ključna elementa, tj. iskano iztočnico (št. 1) in element, ki ga želimo prikazati v obliki kolokatorja (št. 2). Recipročnost razmerja (*DUAL*) pomeni, da če je npr. naša iztočnica pridevnik (št. 2), bodo v skici prikazani vsi samostalniški kolokatorji (št. 1), in obratno, če nas zanima samostalnik, bodo prikazani vsi pridevniški kolokatorji. Recipročnost je upoštevana tudi pri poimenovanju razmerja.

Rezultat tega razmerja je en stolpec besedne skice za posamezni samostalnik ali pridevnik (iskano lemo), ki za samostalniško lemo poda najpogostejše pridevnike, ki jo modificirajo, za pridevniško lemo pa najpogostejše samostalnike, ki ji sledijo. Na sliki 2 so tako na levi navedeni najpogostejši pridevniški kolokatorji za lemo *ogenj*, na desni pa najpogostejši samostalniški kolokatorji za lemo *navzkrižen*, ki so razvrščeni glede na izpostavljenost.

¹ Statistična izpostavljenost je prilagojena oblika Diceovega koeficienta in v grobem temelji na seštevku sopojavljanja besed v nekem slovničnem razmerju. Natančna enačba je dostopna na http://trac.sketchengine.co. uk/attachment/wiki/SkE/DocsIndex/ske-stat.pdf?format=raw. (Dostop 8. 7. 2009.)

ogenj	navzkrižen
Fida PLUS 620m freq = 39083	Fida PLUS 620m freq = 21

a modifier	8961	1.3
taboren	580	96.24
navzkrižen	487	83.4
zmeren	<u>491</u>	68.16
ibek šibek	557	65.35
blag	227	56.11
olimpijski 🗆	460	51.53
topniški 🔳	<u>82</u>	50.9
majhen	784	47.4
peklenski	102	47.04

Slika 2: Vzorec rezultatov za recipročno razmerje pridevnik-samostalnik.

Na dan pisanja tega članka je v slovenski SBS definiranih 18 slovničnih razmerij, osem izmed teh je recipročnih.² Trenutno je v pripravi nova verzija SBS, ki bo obsegala več slovničnih razmerij (vključena bodo recimo razmerja z vezniki), tokrat s slovenskimi poimenovanji in izboljšavami, ki so se pokazale za smiselne pri delu s SkE v začetni fazi gradnje leksikalne baze za slovenščino v okviru projekta Sporazumevanje v slovenskem jeziku.³

Priprava SBS je jezikovnotehnološki postopek, ki ni vezan samo na poznavanje tipičnih jezikovnih pojavov v konkretnem jeziku in nabor oznak, uporabljenih v posameznem korpusu, temveč je pri tem zelo ali celo predvsem pomemben namen, za katerega želimo orodje oz. funkcijo izrabljati. Večina slovničnih razmerij v slovenski SBS, tako trenutni kot nastajajoči, je naravnana izrazito leksikografsko, se pravi za namen izdelave leksikalne baze za slovenščino. Slovnična razmerja so torej definirana glede na vnaprej določene vzorce, za katere predpostavljamo, da so v (slovenskem) jeziku tipični in posledično zanimivi za izdelavo slovarjev in podobnih leksikografskih del.

Zgornji pristop k definiranju slovničnih razmerij seveda ni edini možen. Kot primer

manj vnaprej določujočega pristopa navajamo del skice iz poskusnega 20-milijonskega dela korpusa FidaPLUS za lemo *mesto* (Slika 3).

V ozadju te skice je slovnica, ki je nastala po vzoru slovaške slovnice besednih skic, ki jo je napisal Vladimir Benko in jo je Simon Krek poskusno prilagodil za slovenščino. Primerjava razmerij, opredeljenih na osnovi slovaške SBS, z razmerji iz slovenske SBS, ki je trenutno v rabi za 621-milijonski korpus FidaPLUS, pokaže, da so v slovaškem primeru razmerja zastavljena ohlapnejše. Iskani kolokatorji namreč niso ločeni glede na besedne vrste ali druge slovnične/jezikovne oznake, ampak so v skupni stolpec združeni samo glede na mesto pojavljanja v relaciji do iskane leme.

Na sliki 3 so v levem stolpcu prikazani izrazi, ki se pojavljajo pred povezavo predloga *za* in iskano lemo *mesto*, v desnem pa izrazi, ki

mesto Fida PLUS 20m SLOVASKI freq = 20946

Y za-d X	1543	1.3	v-d Y X	2438 1.3
boj	63	34.65	nov	506 44.85
tekma	81	30.5	središče	145 44.62
dežuren	16	30.32	glaven	78 30.44
potegovati	23	29.4	knežji	10 24.79
kandidat	27	23.8	bližina	28 23.66
napredovati	12	21.05	osvojiti	27 22.4
razpis	18	20.31	okolica	21 20.29
boriti	12	19.32	zasesti	16 20.03
služba	<u>16</u>	15.67	velik	73 19.19
veljati	11	11.9	center	30 18.4
da	12	8.98	večen	10 17.08
iti	8	5.84	francoski	14 16.32
1e	11	5.74	rojsten	<u>10</u> 16.16
v	79	5.49	star	26 15.29
in	18	4.56	svet	<u>8</u> 13.06
se	37	4.07	naš	33 13.0
biti	115	2.82	nemški	13 12.99
še	8	1.1	bližnji	9 12.71
tudi	9	0.44	majhen	15 11.49
			ta	90 10.68
			italijanski	<u>8</u> 10.53
			drug	30 9.91
			nek	11 9.17
			sam	15 8.85
			del	<u>17</u> 8.22

Slika 3: Del skice za lemo *mesto*.

² Več o drugih vrstah razmerij glej v Krek, Kilgarriff 2006.

³ www.slovenscina.eu.

se pojavljajo med predlogom v in iskano lemo *mesto*. Iskana lema je zmeraj označena z oznako X, njena sobesedilna okolica, ki jo želimo prikazati, pa z oznako Y. Ker razmerje ni omejeno na nobeno od besednih vrst, se med njimi pomešano znajdejo samostalniki, glagoli, pridevniki, členki, vezniki, predlogi in zaimki.

Ker se poskuša oddaljiti od obstoječih zakonitosti in predvidevati čim manj, tak pristop po eni strani pušča odprtih več možnosti, da odkrijemo obrobnejše, morda nepričakovane pojave v jeziku, po drugi strani pa je prav zaradi tega potrebno več pozornosti in časa za pregled in interpretacijo množice nediferenciranih jezikovnih podatkov, zaradi česar je v praksi, predvsem v splošnoleksikografskih projektih,⁴ kjer so časovni in finančni okviri jasno zastavljeni, tak pristop manj primeren.

3.2 Označenost korpusa

Za iskanje in povezavo besed v slovnična razmerja, ki so nato predstavljena v besednih

glasben Fida PLUS 620m freq = 149951

post verb:	2563 0.1
spotiti	<u>549</u> 85.35
□ voščiti	<u>405</u> 73.9
□ večeriti	<u>136</u> 62.93
□ zvrstiti	<u>213</u> 54.88
podlagati	<u>61</u> 42.51
skupiniti	<u>69</u> 36.92

skicah, mora SkE vedeti, kako najde besede, ki so v nekem razmerju. Če npr. razmerje predvideva zvezo pridevnika in samostalnika, potrebujemo oblikoskladenjsko označen korpus, kakršen je FidaPLUS (gl. nadaljevanje). Strogo gledano je sicer edini predpogoj za vključitev korpusa v orodje SkE tokeniziranost, tj. razdeljenost besedila na pojavnice, kar pomeni, da bi lahko zelo osnovna oblika SBS delovala že na korpusu, ki ni ne lematiziran ne oblikoskladenjsko označen, seveda pa bi tako dobili zelo osiromašene in leksikografsko ne prav uporabne rezultate.

Ker so besedne skice za slovenščino izdelane na osnovi oblikoskladenjskih oznak, je torej zelo pomembna pravilna označenost korpusa. FidaPLUS je bila označena avtomatsko, kar s seboj prinaša določen odstotek napak, saj je slovenščina zaradi obremenjene morfologije, pri kateri prihaja do prekrivnosti končnic, zelo težaven jezik za označevanje.⁵ Posledica tega je pojavljanje napak v besednih skicah, kjer se lahko v nekem razmerju znajdejo izrazi, ki tja ne sodijo (Slika 4).

Primeri na sliki 4 prikazujejo napačno uvrščene kolokatorje, ki so posledica napačno pripisane leme in oblikoskladenjske oznake, saj je samostalniški kolokator prepoznan kot glagol in je tudi umeščen v stolpec z razmerjem pridevnik-polnopomenski glagol.⁶ Spodnje okno

Dnevnik,časopis	19.20 Kulturni utrip, ponovitev - 19.50	Glasbeni	spoti - 20.00 Pop mix z Andrejo, glasbena
Dnevnik,časopis	živo iz studia Televizije Celje - 21.00	Glasbeni	spoti - 21.30 Oglasi - 21.35 Ko karte govorijo
Dnevnik,časopis	živo iz studia Televizije Celje - 21.35	Glasbeni	spoti - 22.05 Oglasi - 22.10 Viža tedna
časopis	Rock TV 22.00 Legalizacija marihuane 22.37	Glasbeni	spoti 22.40 Poročila: Objektiv Gorenjske
časopis	zdravnik, ponovitev 18.40 Santa Barbara 19.30	Glasbeni	spoti 20.00 Detektiva, 1. del angleške
časopis	vsakogar, ponovitev 18.40 Santa Barbara 19.30	Glasbeni	spoti 20.00 Lutkovno gledališče, prenos
časopis	ponovitev filma 18.00 Avtodrom MMTV 19.30	Glasbeni	spoti 20.00 Očetje in sinovi, ameriški
časopis	prijatelji, ponovitev 19.00Kuhajmo skupaj 19.30	Glasbeni	spoti 20.00 Beverly Hills Buntz, ameriška

Slika 4: Napačna lematizacija kolokatorjev spot, voščilo, večer, zvrst, podlaga in skupina.

- 4 Z izrazom splošnoleksikografski projekti so mišljeni projekti, usmerjeni v opis splošnega jezika.
- 5 O tem, da je slovenščina v primerjavi z angleščino težaven jezik za označevanje, priča tudi dejstvo, da slovenski nabor oblikoskladenjskih oznak obsega okoli 2.000 elementov (Erjavec, Krek 2008: 51). Za primerjavo, standardni nabor britanskega nacionalnega korpusa vsebuje okoli 160 oznak (Označevalni sistem CLAWS).
- 6 »Spoti« je mogoče razumeti kot samostalnik v množini in kot obliko glagola spotiti.

s konkordancami potrjuje, da gre v resnici za razmerje *pridevnik-samostalnik*, k napaki pa je najverjetneje prispevalo tudi pomanjkanje konteksta, ki pomaga pri avtomatskem razdvo-umljanju.

Naslednji izsek besedne skice za glagol pobirati (Slika 5) prikazuje razmerje samostalnik v imenovalniku-polnopomenski glagol, v katerem želimo najti osebke, ki se tipično pojavljajo z zadevnim glagolom. Na vrhu zadetkov smo označili kolokatorja, ki sta sicer pravilno lematizirana, vendar napačno označena kot samostalnika v imenovalniku namesto kot samostalnika v tožilniku, kar potrdijo konkordance v spodnjem oknu.

Nekatera mesta s pogosto napačno označenimi besedami so:

- a) pridevnik ženskega spola v tožilniku v krepkem tisku je napačno označen kot prislov (... označili ... za »pogumno žensko in veliko državnico.«);
- b) samostalnik srednjega spola je napačno označen kot glagolska oblika v tretji osebi srednjega spola (*Kakovost slike je osnovno merilo pri nakupu digitalnega fotoaparata.*);

pobirati Fida PLUS 620m freq = 11487

<u>has subj</u>	<u>1675</u> 7.6
denar	<u>69</u> 26.29
☐ država	<u>106</u> 25.51
☐ davek	<u>34</u> 24.55
□ suša	<u>16</u> 23.9
☐ davkarija	<u>9</u> 23.82

c) samostalnik v imenovalniku ali tožilniku množine je napačno označen kot samostalnik v rodilniku ednine (*Nekje sem zasledila, da ženske že po naravi težje prenašamo vročino ...*; *Dobro poznam ta občutek: gledaš mlade ženske, toda one te sploh ne opazijo.);*

- d) samostalnik moškega spola v tožilniku je napačno označen kot samostalnik moškega spola v imenovalniku (*Prispevek pobira javno komunalno podjetje* ...);
- e) pridevnik moškega spola v imenovalniku je napačno označen kot samostalnik srednjega spola v rodilniku množine (*Pričakovali bi, da bo sin vesel takšne materine požrtvovalnosti.*).

Orodje SkE sicer prepoznava tudi odvisnostnoskladenjsko razčlenjeni korpus, za kakršnega je mogoče napisati še kompleksnejšo in s tem natančnejšo SBS od trenutno veljavne za slovenščino, ki upošteva le oblikoskladenjske oznake (Krek, Kilgarriff 2006). V tem primeru so tudi rezultati, dobljeni z besedno skico, pravilnejši, saj so prav odvisnostnoskladenjska razmerja ta, ki bi jih v besedne skice radi zajeli, vendar trenutno nimamo dovolj natančnega sredstva za njihovo razločevanje. Če bi lahko odvisnostnoskladenjske oznake uporabili vzporedno z oblikoskladenjskimi, bi to zelo zmanjšalo število napak zaradi napačne lematizacije.

Delo,časopis	toliko kot pri nogometnem prvenstvu. Denar pobira	Karmen. Naj zmaga najboljši, C s se glasi
revija	lahko » z lahkoto in z namenom « sploh še pobira	javni denar za neko poslanstvo, ki ga preprosto
revija revija	organizirane kriminalne skupine, ki revežem pobirajo	denar z obljubami, da jih bodo spravile
revija	tolarjev. < Navaden > Denar sicer pobira	žalsko komunalno podjetje, vendar je strogo
Dnevnik,časopis	www.dnevnik.si. (bš Kdo staršem pobira	denar? Osupel sem v Dnevniku (8.

Slika 5: Napačna prepoznava sklonov kolokatorjev.

7 Za korpus FidaPLUS, kot je v SkE, velja nabor oznak MULTEXT-East v.3, v prihodnje pa bo uporabljen izboljšani nabor oznak JOS v kombinaciji z natančnejšim označevalnikom (Erjavec, Krek 2008).

3.3 Besedilna sestava korpusa

Ena temeljnih reči, ki najbolj vpliva na to, kakšne rezultate bo dala besedna skica, je besedilna sestava korpusa, saj analiza ne more pokazati ničesar drugega kot preoblikovano sliko tega, kar je bilo vanjo poslano. Pod besedilno sestavo razumemo vrste besedil, ločene glede na zvrst (leposlovje, stvarna besedila) in prenosnik (govorno, časopis, internet itn.), njihov obseg ter razpršenost virov znotraj posamezne vrste besedil.

FidaPLUS kot referenčni korpus vsebuje zelo raznolika besedila, med njimi tudi televizijske sporede, male oglase ipd. Tovrstna besedila zaradi svoje ponavljajoče se narave vsebujejo veliko istovrstnih zadetkov, poleg tega pa v njih tudi primanjkuje konteksta, ki bi pomagal razdvoumiti, zaradi česar so včasih rezultati, dobljeni iz takšnih besedil, manj relevantni za splošnoleksikografske namene. Danes se tovrstna besedila (vključno s cenami delnic, križankami, športnimi rezultati, kazali ipd.) vsaj pri leksikografskih korpusih običajno odstrani (Atkins, Rundell 2008: 85, Kilgarriff, Rundell, Dhonnchadha 2006).

Na sliki 6 so označeni kolokatorji, ki, kot potrjujejo konkordance, izvirajo predvsem iz omenjenih besedil.

glasben Fida PLUS 620m freq = 149951

161922 70 -1 --- 16-1 2422 21

modifies	<u>161833</u> 7.0	adj modified	3432 2.1	
☐ lestvica	<u>12610</u> 86.02	□ zabavno	<u>217</u> 69.37	
oddaja oddaja	<u>9747</u> 70.51	☐ slovensko	<u>201</u> 56.87	
poslušalnica	339 67.28	☐ srednje	<u>174</u> 53.25	
□ jutranjica	<u>349</u> 66.34	☐ domače	<u>129</u> 51.8	
motorček.	<u>485</u> 64.71	☐ literarno	<u>48</u> 48.72	
□ šola	<u>15271</u> 61.99	☐ samostojno	<u>60</u> 42.25	
utrineSlovenija	<u>152</u> 58.56	☐ humoristično	<u>17</u> 41.14	
	128 56.58	☐ ljudsko	<u>25</u> 39.17	
utrinekSlovenija	120 30.38	□ nizko	<u>54</u> 38.99	
utrinek	<u>1063</u> 56.29	□ veliko	269 38.56	
spremljava 🗆	<u>807</u> 55.21		10 27 00	
□ scena	<u>2255</u> 55.02	narodnozabavno	<u>12</u> 37.98	
☐ video	<u>1552</u> 54.33	☐ razvedrilno	<u>12</u> 36.27	
□ zvrst	<u>1000</u> 54.07	nekaj	<u>120</u> 36.25	
□ želja	<u>3733</u> 53.04	☐ bogato	<u>47</u> 35.98	
pedagog	<u>651</u> 51.81	popularno popularno	<u>20</u> 35.94	
cocktail	<u>154</u> 51.48	osnovno	<u>27</u> 34.57	
		1100 771 1	1. 1	12.20 (1) 1
Dnevnik,časopis		-		13.30 Glasbeni motorček, ponovitev mladinske oddaje -
Dnevnik,časopis Dnevnik,časopis		•		19.05 Glasbeni motorček, mladinska oddaja - 19.30 Ris 13.30 Glasbeni motorček, ponovitev mladinske oddaje -
Dnevnik,časopis Dnevnik,časopis	-	-		19.05 Glasbeni motorček, mladinska oddaja - 19.30 Ris
Onevnik,časopis				13.30 Glasbeni motorček, ponovitev mladinske oddaje -
Dnevnik,časopis	-	-		19.05 Glasbeni motorček, mladinska oddaja - 19.30 Ris
časopis	22.40 Halifax,	nani. T	V 3	14.00 Glasbeni motorček - 16.30 Glasba brez meja - 17
časopis	Popotov	anja z Janinom - 18	.00 Bonanca -	19.05 Glasbeni motorček - 19.30 Risanke - 20.00 Posle

Slika 6: Prikaz tipičnih kolokatorjev iz televizijskih sporedov.

časopis

Južni Brooklyn. TV 3 13.30 Glasbeni motorček - 17.00 Poletni utrip - 18.00

Poleg zgoraj omenjenega na dobljeno skico odločilno vpliva tudi odločitev snovalcev korpusa glede velikosti posameznih besedilnih komponent v korpusu (npr. knjižni proti časopisnemu prenosniku). S ponazoritvijo na sliki 7 želimo pokazati, kako se s spremembo besedilne sestave korpusa spremeni tudi besedna skica.

glasben

modifies	30167 3	3.9
☐ lestvica	7648 10	0.06
□ scena	<u>640</u> 6	5.62
□ šola	<u>1765</u> 6	5.32
■ industrija	<u>597</u> 6	5.22
■ zvrst	320	5.9
■ datoteka	<u>363</u> 5	5.86
skupina	950 5	5.55
kariera	280	5.4
plošča	321 5	5.33
spremljava	191	5.2
■ festival	311 5	5.15
založba	227 5	5.15
■ okus	192 4	1.86
Podkorpus rev	rij	

Slika 7: Različni kolokatorji glede na podkorpusa revij in časopisov.

Prikazan je vzorec prvih trinajstih zadetkov za razmerje pridevnik-samostalnik za lemo glasben v dveh podkorpusih: s črno so označeni tisti kolokatorji, ki jih ne najdemo v drugem podkorpusu. Situacija je v tem primeru seveda skrajna, saj imamo na eni strani le podkorpus vseh časopisnih besedil iz korpusa FidaPLUS, na drugi strani pa vsa revijalna besedila iz korpusa, referenčni korpus pa je kombinacija obojega. Kljub temu je razvidno, da z izdelavo dveh podkorpusov oz. tudi z manj izrazito spremembo deležev v referenčnem korpusu sprožimo spremembe v razporeditvi kolokatorjev v besednih skicah. Tako lahko glede na sestavo korpusa FidaPLUS - 65 odstotkov besedil nosi oznako časopisno, 23 odstotkov oznako revijalno -8 sklepamo, da bi se s krčenjem časopisnega dela korpusa zadetki iz televizijskih sporedov znašli *nižje* na seznamu kolokatorjev.

Televizijski sporedi se namreč pogosteje pojavljajo v časopisnem kot pa v revijalnem gradivu, o čemer pričajo podatki o izvoru levo od konkordanc na sliki 6.

3.3.1 Razpršenost virov

Kot že omenjeno, je za to, da dobimo relevantne podatke, pomembno tudi, da so viri čim bolj razpršeni. Besedna skica ne nudi podatkov o razpršenosti virov, v katerih se pojavljajo kolokatorji, je pa ta podatek mogoče najti v konkordancah.

Za kolokator *kredit* (Slika 8), ki je del razmerja *samostalnik-predlog 'glede'*, ob vpogledu v konkordance ugotovimo, da vsi zadetki izhajajo iz enega vira. Kadar so viri za zadetek nerazpršeni, kolokator za splošnoleksikografske namene najverjetneje ni relevanten.

3.4 Statistično razvrščanje zadetkov v besedni skici

Poleg besedilne sestave in označenosti korpusa, na kar kot uporabniki nimamo vpliva, je za rezultate, ki jih dobimo v besedni skici, pomembno tudi to, kakšne nastavitve izberemo pri uporabi orodja. Vsakič ko želimo izdelati besedno skico, izberemo tudi možnosti prikaza, kar presega zgolj število izpisanih kolokatorjev ipd. (Gantar idr. 2009: 34), ampak zadeva tudi določitev tega, kolikšni naj bosta vrednosti najnižje absolutne pogostosti kolokatorjev in statistične izpostavljenosti oz. po katerem od statističnih izračunov naj bodo kolokatorji v skici razvrščeni.

Na sliki 9 prikazujemo mesta, kjer se skici za isto iskano lemo razlikujeta glede na izbrani statistični izračun.

Z besednih skic za lemo *obdobje*, okrajšanih na 16 kolokatorjev, je v slovničnih razmerjih *samostalnik-samostalnik v rodilniku* in *pridevnik-samostalnik* mogoče razbrati, da se skici precej razlikujeta – s črnimi kvadratki so

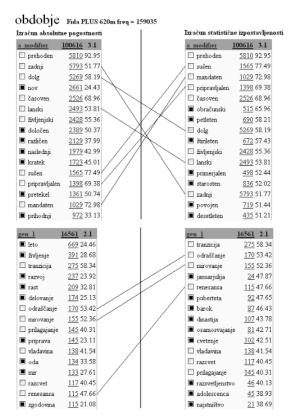
⁸ V novem referenčnem korpusu slovenščine, ki nastaja v okviru projekta Sporazumevanje v slovenskem jeziku, so besedilna razmerja zastavljena drugače, predvsem je načrtovan manjši obseg časopisne komponente (Logar Berginc, Šuster 2009).

Vprašanje Fida PLUS 620m freq = 331775

a modifier	123348	2.2	post glede-d	1006 25.6	post o-d	<u>10554</u> 9	9.9
nerešen nerešen	<u>1554</u>	72.33	kredit	43 35.99	smiselnost	<u>155</u> 55	5.67
referendumski	1737	70.62	usposabljanje	<u>17</u> 23.75	ueroizpoved veroizpoved	134 54	4.86
☐ pogost	3672	70.23	☐ meja	<u>19</u> 16.37	☐ smisel	<u>148</u> 33	3.45
□ odprt	5049	69.51	☐ financiranje	<u>9</u> 14.48	umestnost	13 29	9.86
☐ zastavljen	1708	69.32	☐ zdravje	<u>10</u> 13.83	primemost	<u>32</u> 29	3.35
mandaten	827	66.59	pravilo	<u>9</u> 12.42	upravičenost	<u>36</u> 29	9.08
nagraden n	2464	64.56		>>	dopustnost	<u>15</u> 28	3.98
				_			

FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef Kvas Odgovor župana na svetniško vprašanje glede stanovanjskih kreditov Ob dnevu samostojnosti
FPsplet_internet Jožef

Slika 8: Prikaz kolokatorja in konkordanc iz istovrstnih virov.



Slika 9: Primerjava rezultatov glede na dve vrsti statističnega izračuna za lemo *obdobje* v dveh slovničnih razmerjih.

označeni kolokatorji, ki jih skica pri drugem statističnem izračunu ne pokaže, s črtami so povezani kolokatorji, pri katerih je porazdeljenost precej različna. Beseda, ki je v korpusu relativno redka, a se pogosto pojavlja prav z iskano lemo, bo imela zelo visoko vrednost statistične izpostavljenosti. Tako se na sliki v drugem stolpcu znajde kolokacija obračunsko obdobje, ki se v korpusu pojavi 515-krat, kar je premalo, da bi se uvrstila na vrh seznama zgolj po absolutni frekvenci (levi stolpec). Sklenemo lahko, da pri prvem izračunu najdemo pri vrhu bolj splošne kolokatorje (zadnji, dolg, nov), pri drugem pa kolokatorje, ki so sami zase redkejši v korpusu (sušen, mandaten, pripravljalen) (gl. tudi Atkins, Rundell 2008: 110). Iz vsega tega je mogoče zaključiti, da je izbira statističnega izračuna in ostalih nastavitev, enako kot večina drugih stvari, prav tako odvisna od namena uporabe.

4 Zaključek

V prispevku smo predstavili možnosti, ki jih ponuja uporaba besednih skic pri raziskovanju jezika, s poudarkom na leksikografskem delu. Splošni način njihove uporabe in zlasti njihova neprecenljiva vrednost pri leksikografskem delu sta bila v preteklosti že opisana, zato smo v tem prispevku predvsem opozorili na uporabniku manj očitne vidike in ozadje njihovega delovanja, ki pa pomembno vplivajo na dobljene rezultate. To na prvi stopnji obsega besedilno sestavo in označenost korpusa ter razpršenost virov, ki morajo biti dobro premišljeni, da odsevajo čim bolj realno podobo jezika. Dobljeni rezultati so nadalje odvisni od načina obdelave podatkov, kjer igra ključno vlogo slovnica besednih skic, in čisto na koncu tudi od uporabnika samega, ki na rezultate vpliva z izbiro nastavitev orodja. Neizpodbitno dejstvo je namreč še vedno to, da kljub neizmernemu delu, ki ga z uporabo tega in podobnih orodij za nas opravi računalnik, za vnosom podatkov, potrebnih za nastanek besednih skic in jezikovnotehnoloških aplikacij nasploh, še

zmeraj stoji človek, kar pomeni, da to, kar surovim oziroma neobdelanim jezikovnim podatkom pripišemo, tudi (v nekoliko spremenjeni obliki) iz njih dobimo nazaj.

Literatura in spletni viri

- ATKINS, B. T. Sue, RUNDELL, Michael, 2008: *The Oxford Guide to Practical Lexicography*. Oxford: Oxford University Press.
- ERJAVEC, Tomaž, KREK, Simon, 2008: Oblikoskladenjske specifikacije in označeni korpusi JOS. Erjavec, Tomaž, Gros, Jerneja (ur.): Konferenca Jezikovne tehnologije. Ljubljana: Institut Jožef Stefan. 49-53.
- GANTAR, Polona idr., 2009: Specifikacije za izdelavo leksikalne baze za slovenščino: opis analize referenčnega korpusa: http://www.slovenscina.eu/Vsebine/Sl/Kazalniki/Kazalnik5/Kazalnik5.pdf. (Dostop 8. 7. 2009.)
- KILGARRIFF, Adam, RUNDELL, Michael, UÍ DHONNCHADHA, Elaine, 2006: Efficient corpus development for lexicography: building the New Corpus for Ireland. Language Resources and Evaluation Journal 40/2. 127-152.
- KREK, Simon, KILGARRIFF, Adam, 2006: Slovene word sketches. Erjavec, Tomaž, Gros, Jerneja (ur.): Konferenca Jezikovne tehnologije. Ljubljana: Institut Jožef Stefan. 62-67.
- LOGAR BERGINC, Nataša, ŠUSTER, Simon (v tisku): Gradnja novega korpusa slovenščine. Jezik in slovstvo 54: 3-4.
- Označevalni sistem CLAWS: http://ucrel.lancs.ac.uk/claws/. (Dostop 3. 7. 2009.)
- Sketch Engine beta: http://beta.sketchengine.co.uk/. (Dostop 3. 7. 2009.)