

Milica Gačić

KORIŠTENJE KONKORDANCIJA U IZRADI PEDAGOŠKIH RJEČNIKA STRUKE

Stručni članak. Primitljen 1. 2. 1985.

UDK 37:800.866(038):801.82:681.3

Za efikasniju nastavu jezika struke nužna je izrada materijala s kojima će se postići optimalno učenje jezika. Korištenje kompjutera može nam znatno olakšati izradu pedagoških rječnika struke — bilo u klasičnoj formi standardnog rječnika — s tim da su obuhvaćene najrelevantnije i najčešće natuknice koje se koriste u određenoj struci, ili u obliku kontekstualnog rječnika struke koji može pružiti izuzetno vrijedne primjere upotrebe leksičke jedinice, kolokacije ili idioma.

Korištenje konkordancija je jedan od metodoloških postupaka koji može pridonijeti rješavanju bitnih pitanja u vezi s izradom rječnika struke jer pružaju kontekstualne podatke o riječima.

Prva dekada automatske izrade konkordancija teksta su pedesete godine ovoga stoljeća. Među pionirima u korištenju te metodologije svakako je Roberto Busa, teolog koji je proučavao radove Tome Aquinskoga. Njegovim radom postavljeni su temelji pa je 1956. godine u mjestu Gallarate ustanovljen Centar za automatizaciju literarnih analiza (Centro per L'Automazione dell' Analisi Letteraria).¹ U tom centru 1964. godine radilo je šezdeset ljudi.

U radu *Nieuwe wegen in der lexicologie* (1963), holandski autor F. de Tollenaere naglašava da su leksikografija, izrada indeksa riječi i izrada konkordancije tri aktivnosti što pripadaju u *leksikologiju* koju definira kao »teoriju rječnika koja uključuje semantiku i morfologiju, tvorbu riječi, etimologiju, sinonimiju itd.«

Pod pojmom pedagoškog rječnika razumijevamo funkcionalni profil rječnika određene struke (u pravilu dvojezični) koji obuhvaća leksičku građu što optimalno zadovoljava potrebe nastave i svladavanja jezika struke. To znači da takav rječnik treba da sadrži stupnju nastave stranog jezika primjeren broj leksičkih jedinica optimalne korisnosti s obzirom na mogućnost razumijevanja stručne informacije na stranom jeziku. Time je pedagoški rječnik prije svega pomoć nastavniku zbog ograničenog dometa empirije i intuicije, te studentu zbog optimalizacije uloženog truda i sredstava u učenju stranog jezika.

Više je načina prikupljanja rječničke građe. Kad je riječ o problematici rječničke građe za potrebe nastave jezika struke, izuzetnu nam pomoć može pružiti kompjutorska obrada reprezentativnog korpusa jezika struke. U uvje-

¹ Burton, D. M., str. 12.

tima kad za određenu struku postoji jezični korpus unesen na magnetsku traku (ili drugi medij), relevantne podatke za izradu pedagoškog rječnika struke možemo dobiti ako kompjutorskom obradom dobijemo ove ispise:

- (1) abecednu listu riječi,
- (2) listu riječi po frekvenciji,
- (3) konkordanciju teksta.

U sklopu raspoloživog vremena i u skladu s psiho-fizičkim mogućnostima studenata, kroz 120 sati nastave možemo očekivati usvajanje dosta ograničenog broja riječi. Da bismo utvrdili leksički okvir unutar kojega se treba kretati da bi se pažnja posvetila komunikacijski važnijim dijelovima leksičke građe, kao jedan od osnovnih kriterija može nam poslužiti frekvencijska lista leksičkih jedinica.

Korisnost utvrđivanja frekvencija jezičnih sredstava uočena je i korištena još u starom vijeku, a prva brojanja riječi rađena su za potrebe slijepih i stenografije. Za pedagoške potrebe prva brojanja riječi obavljena su u prvoj polovici 20. stoljeća.

Frekvencijski podaci o riječima u općem jeziku koristili su se u nizu tečajeva i projekata za učenje jezika. No vrijednost takvih podataka mnogo je upotrebljivija za pedagoške potrebe u jeziku struke jer jezik struke ima ograničen leksički repertoar koji je moguće dosta pouzdano utvrditi korištenjem dobrog korpusa.

Ovdje ćemo navesti samo neke podatke koje ističu razni autori u vezi s korištenjem frekvencija. Tako P. Guiraud² kaže:

- 100 prvih riječi pokriva 60% teksta,
- 1 000 prvih riječi pokriva 85% teksta,
- 4 000 prvih riječi pokriva 97,5% teksta.

V. Pravda³ navodi da minimalni vokabular od približno 2 200 riječi pokriva, prema provedenim ispitivanjima, oko 95% teksta, a u priručniku *Languages for Special Purposes*⁴ autori navode da lista od 2 000 općih riječi i 500 specijaliziranih termina pokriva oko 90% stručnog teksta.

U vezi s dosad provedenim frekvencijskim ispitivanjima jezika struke treba naglasiti da je minimalne vokabulare u jeziku struke lakše utvrditi na osnovi podataka o frekvenciji jer, s obzirom na uži izbor leksičkih sredstava, postoje manje razlike u distribuciji.

Osvrnut ćemo se na primjer jedne takve analize.⁵ Analizirajući frekvencijsku listu korpusa krivičnih disciplina utvrdili smo da bi pedagoški bilo opravdano obraditi leksičke jedinice do minimalne vrijednosti frekvencije $f = 4$ u korpusu od 120 000 pojava. Za tu smo se vrijednost odlučili zato što smatramo da taj broj javljanja neke riječi u korpusu nije slučajnost. Tako dobivamo određeni broj leksičkih jedinica koje se mogu svrstati u tri kategorije:

- (a) opće riječi,
- (b) opće riječi u stručnoj upotrebi,
- (c) stručne riječi.

² Guiraud, P., str. 10.

³ Pravda, V., str. 270.

⁴ *Languages for Special Purposes*, str. 34.

⁵ Gačić, M., *Frekvencijsko-leksička analiza engleskog kao jezika struke krivičnih disciplina*, neobjavljena doktorska disertacija, Zagreb, 1982.

Prvu kategoriju čine uglavnom riječi obuhvaćene prijašnjim stupnjevima učenja, pa u okviru nastave jezika struke nije predviđen proces njihovog aktivnog usvajanja, već samo ponavljanja, no one treba da budu uvrštene u pedagoški riječnik. Druge dvije kategorije mogu se obuhvatiti nastavom jezika struke, odnosno nužno im je posvetiti pažnju da bi bile aktivno usvojene.

Osnovni lingvistički postupci kojima treba podvrći tako relativno mehanički određen frekvencijski popis jest *lematizacija* (odnosno postupak svođenja kosih oblika na kanonski rječnički oblik — jedinica za imenice, infinitiv za glagole, pozitiv za pridjeve).

Tako dobivene rječničke natuknice treba svrstati abecednim redom. Broj leksičkih jedinica u lematiziranoj varijanti smanjuju je otprilike 26% u engleskom, a prema nekim podacima 20% u njemačkom jeziku.⁶ U našem slučaju (u krivičnim disciplinama) broj leksičkih jedinica sveden je sa 3 379 na 2 494 jedinica.

Abecedni popis može sada, osim kao osnova za rječničke natuknice, poslužiti za uspješnije pronalaženje odgovarajuće leksičke jedinice u konkordanciji.

Ž. Bujas⁷ ovako definira konkordanciju: »Konkordanca je tako priređen tekst da je svaka pojedina riječ tog teksta (bez obzira na to koliko se puta ponavljala, gdje se javljala i u kojem obliku dolazila) prikazana abecednim redom sa svojim kontekstom. Veličina konteksta u konkordancama je proizvoljna, ali — premda kontekst može obuhvaćati cijelu strofu, rečenicu ili odlomak — obično je u opsegu jednog stiha ili jednog retka kompjuterske liste (8 do 10 riječi ispred i iza konkordirane riječi)«.

Za potrebe analize jezika struke krivičnih disciplina izrađena je potpuna tzv. normalna kontekstualna konkordancija. Ispis sadrži ove podatke za svaku riječ: oznaku za knjigu, broj stranice u knjizi, broj sloga u ispisu teksta (red), tekst koji prethodi, razmakom istaknutu riječ koju konkordiramo i tekst koji slijedi.

Na sljedećoj stranici dajemo primjer opisane konkordancije.

⁶ Teubert, W., str. 296.

⁷ Bujas, Ž., *Kompjuterska konkordanca . . .*, str. 35.

K O N K O R D A N C A T E K S T A

KNJC STR RED	TEKST KOJI PRETHODI	RIJEČ	TEKST KOJI SLJEDI
K9 17 14	O THIS END, OFFICERS USING CYCLES, MOTOR-CYCLES AND VANS	CARRY OUT SPOT CHECKS AROUND THE BANKS. IF THE HEADS	
K9 45 15	T ONLY IN CONNECTION WITH ACTIVE POLICE WORK BUT ALSO TO	CARRY OUT TASKS OF AN ADMINISTRATIVE AND FINANCIAL N	
K9 53 04	ENQUIRIES ALONG THE RIGHT LINES. IT IS ALSO POSSIBLE TO	CARRY OUT SEARCHES FOR ACCOMPLICES, I.E. TO COMPIL	
K3 580 06	MILITARY TRAINING AND DISPLAY, THE WIDESPREAD PRACTICE OF	CARRYING GUNS OR KNIVES; THE CONCERN WITH PERSONAL H	
K3 586 06	IMPROVING LOCKING DEVICES OR DOING AWAY WITH THE NEED OF	CARRYING CASH. THESE IDEAS DO NOT REPRESENT AN ABRUP	
K6 169 02	THIS NATURE - IN ORDER TO REDUCE THE VISIBILITY OF THOSE	CARRYING A HYPODERMIC NEEDLE BECAUSE HE OR SHE IS DI	
K8 32 21	ALLERGY, E.G. SHIFFLES, RUNNING EYES; OR A PERSON MAY BE	CARRYING THEM, MUST BE REMOVED TO THE LABORATORY FOR	
K8 32 23	SCRIEBE IN CHAPTER 15. SUSPICIOUS STAINS, OR THE OBJECTS	CARRYING BLOOD, AND OTHER SPORTS THAT MAY COLLECT AN	
K8 120 03	FF AND PLACED IN A GLASS VIAL. CRACKS IN THE FLOOR, SOIL	CARRYING THE IMPRESSION, NOTING THE LOCATION ON THE	
K3 372 03	ION CAN BE AVOIDED BY CAREFUL EXAMINATION OF THE SURFACE	CARS STOLEN THE KEY HAS BEEN LEFT IN THE IGNITION.)	
K3 384 02	AT A LARGE SHARE OF CITIZEN COMPLAINTS RADIOED TO PATROL	CARS WITH MISMATCHED HUB CAPS, OR DIRTY CAR WITH CLE	
K3 600 12	TO AVOID DIRECT EYE CONTACT PERSON WEARING COAT ON HOT DAYS	CARS, WELL-LIGHTED STREETS, ALARM SYSTEMS, AND PROP	
K3 600 19	ON SCREENS; PROVIDING MORE EFFECTIVE LOCKING DEVICES FOR	CARS, HOUSES, AND OFFICES; SURROUNDING RESIDENTIAL A	
K4 44 04	S AND EXPLOSIVES. 12. RECOGNITION AND RECOVERY OF STOLEN	CARS, 13. TRAFFIC REGULATION AND CONTROL. 14. PURSUI	
K6 108 20	CATION BETWEEN POLICE HEADQUARTERS AND THE ROYING PATROL	CARS IS PROBABLY THE MOST IMPORTANT OF THEM ALL. THE	
K6 108 22	ITY OF THOSE CARRYING OUT THE CHECKS, AND THE DETAILS OF	CARS FOUND INSECURE WERE RECORDED; FROM THIS INFORMA	
K8 410 10	HOW IT WAS POSSIBLE TO IDENTIFY WHAT PROPORTION OF THOSE	CARS ON THE ROAD. THUS, IT MUST BE ACKNOWLEDGED THAT	
K8 466 10	MS OF PASSENGER MILES, OR EVEN IN TERMS OF THE NUMBER OF	CARS OCCASIONALLY EXPLODE WHILE THEY ARE BEING PAINT	
K8 466 10	STRIAL ACTIVITIES, LAPSE STORAGE TANKS AND RAILROAD TANK	CARTONS'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K8 26 10	OWN AREA, SOMETHING WHICH I ARRANGED TO DO WITH SHERIFF	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K7 61 12	BE IN THE WIT, PARTICULARLY IN OUTDOOR CRIMES. CARDGARD	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K7 61 12	OSIBLE TO DETECT THEM FROM THE LOCATION OF THE RECOVERED	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K7 69 02	ES. BOTH SPECIMENS FACE THE SAME DIRECTION. ONE BULLET OR	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K7 69 04	NING PROMINENT STRIATIONS, FOLLOWING THE OTHER BULLET OR	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K8 382 03	A REPEATING MAGAZINE, FOLLOWING THE DISCHARGE OF ONE	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K8 382 03	OF ONE CARTRIDGE, WILL AUTOMATICALLY EJECT THE EXPANDED	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K8 382 03	OPTICALLY EJECT THE EXPANDED CARTRIDGE CASE. LOAD A NEW	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K8 382 12	AND MARKS THE TRIGGER, THE EXTRACTOR, WHICH AT THOMAS THE	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K8 382 12	THE CHAMBER; THE EJECTOR, WHICH AT IDENT AND MARK THE	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K8 382 14	VAL OF THE CASE, A PISTOL MAY LEAVE ADDITIONAL MARKS ON	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K8 394 12	NGE CARTRIDGE, WHICH SERVES AS A RETAINER INTO WHICH THE	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K8 402 07	DIFFERENT FROM THOSE OF MOST OTHER MARKS ON THE SHELL OR	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K8 402 07	REPEATING AND AUTOMATIC WEAPONS ARE DESIGNED TO REMOVE A	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K8 402 09	O THE LATTER, WHILE THESE OPERATIONS ARE PROCEEDING, THE	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K7 61 17	T CASIS FOR COMPARISON IDENTIFICATION. TO SUMMARIZE: THE	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K7 61 17	HATED ARE CHARACTERISTIC FIGURES OF FIRING PINS ON RIFLE	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K7 50 09	S USED FOR RIFLE ONLY FIGURE 4-6; COMPONENT PARTS OF A	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K7 69 01	ACTIVELY AS IT WAS AT THE TIME OF THE CRIME. EACH BULLET OR	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K1 3 22	R PROCEEDINGS, CIVIL OR CRIMINAL, IN RESPECT TO THE SAME	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K1 3 22	S AGAINST THE PERSON ACT 1961, SECTIONS 44, 45). IN THIS	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K1 9 12	POLICY OF THE CRIMINAL LAW TO PREVENT. FOR EXAMPLE, IN A	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K1 15 24	VEHICLE THE ACCUSED WAS UNAWARE OF THIS OCCURRENCE. THIS	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K1 15 24	FECT THE COMPANY WITH LIABILITY, A STILL MORE REMARKABLE	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K1 10 16	S VICARIOUS LIABILITY, TWO OF THE FIVE LAW LORDS IN THIS	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	
K1 24 12	TS; ONLY PARLIAMENT CAN ASSESS THEM. ALTHOUGH M'NAGHTEN'S	CARTON'S SPONSORSHIP. FOR FOUR MONTHS, FOR FOUR HOUR	

U primjeru vidimo da se riječ *cartridge(s)* javlja 20 puta. Da se kolokacija *cartridge case* javlja osam puta, da se izraz *cartrige or shell* (odnosno obrnuto) javlja četiri puta. Naravno, da bismo kolokacije unijeli u pedagoški rječnik, potrebno je najprije definirati što je to kolokacija, na čemu se ovdje nećemo zadržavati.

Konkordancija će nam pružiti nezaobilazne pokazatelje za utvrđivanje značenja riječi (što je riječ »stručnija«, to je jednoznačnija), a posebno kad je u pitanju razdvajanje homografa koji su izrazito česti u općem engleskom jeziku, a nešto manje česti u stručnom engleskom jeziku. Zatim, konkordancija je nezaobilazan izvor za utvrđivanje kolokacije i idioma u kojima se javljaju riječi koje smo izabrali.

Danas je već moguća automatska lematizacija korpusa, npr. u Institutu für deutsche Sprache u Mannheimu razrađen je paket programa Lemma koji je 90% algoritam, što znači da ne pokriva svaku nepravilnost njemačkog jezika.

Spomenuti kompjutorski ispisi omogućuju također stvaranje ekscerpta za rječničke kartice, i tako dolazimo do pitanja što treba da sadrži rječnička natuknica u pedagoškom rječniku. Smatramo da treba da sadrži ove informacije:

- o ortografiji (upozoriti na varijante spelovanja, na primjer)
- o transkripciji (izgovoru)
- o funkciji (morfološkoj i sintaktičkoj)
- o frekvenciji
- o značenju (semantičku)
- o najčešćim kolokacijama, složenicama i idiomima.

Možemo sažeti da je osnovni proces dobivanja pedagoškog rječnika korištenjem automatski izrađenih konkordancija ovakav:

I — unos teksta

- izolacija riječi pojavnica u tekstu
- njihovo rangiranje po frekvenciji
- lematizacija
- abecedno sređivanje
- provjera u rječniku
- provjera u konkordanciji.

II Ispis natuknice:

- obrada (izgovor, vrsta riječi)
- izdvajanje homografa
- traženje kolokacija
- obrada na drugi jezik
- provjera
- konačna verzija.

Vrijedno je još spomenuti da je jedna od mogućih varijanti pedagoškog rječnika za izradu kojeg nam je vrlo korisna konkordancija tzv. kontekstualni rječnik.

Autori kontekstualnog rječnika geologije⁸ oslonili su se na korpus od 40 000 riječi te na osnovi njega sastavili rječnik za uvođenje u geološke studije od 1 500 leksičkih jedinica uz sintaktičku i semantičku analizu.

Shodno iznesenim stavovima zalažemo se za kompjutorsku obradu lingvističkih podataka samo kao za jedan od načina prikupljanja objektivnih podataka (npr. o frekvenciji, distribuciji i rangu riječi) koji može bitno pomoći lingvistima da, koristeći se podacima do kojih nije moguće doći, intuicijom dadu doprinos procesu učenja jezika.

U ovom slučaju to se odnosi na unapređivanje leksikografskog rada u domeni svladavanja stranog jezika struke.

LITERATURA

- Bujas, Ž., »Concordancing as a Method in Contrastive Analysis.« *SRAZ*, 23, str. 49—62, 1967.
- Bujas, Ž., »Kompjuterska konkordanca Gundulićeva »Osmana«, *Filologija*, 7, str. 35—59, 1973.
- Burton, D. M., »Automated Concordances and Word Indexes: the Early Sixties and the Early Centres.« *Computer and the Humanities*, 15, str. 83—100, 1981.
- Burton, D. M., »Automated Concordances and Word Indexes: the Fifties.« *Computer and the Humanities*, 15, str. 1—14, 1981.
- Burton, D. M., »Automated Concordances and Word Indexes: the Process, the Programs, and the Products.« *Computer and the Humanities*, 15, str. 139—154, 1981.
- Centre for Information on Language Teaching and Research, *Languages for Special Purposes*, London, CILT for British Association for Applied Linguistics, 1971.
- Descamps, J. L. i dr., *Dictionnaire contextuel de français pour la géologie*. Paris, Didier, 1976.
- Gačić, M., »Primjena kompjutora u pripremi materijala za analizu jezika struke.« *Zbornik II znanstvenega srećanja računalniška obdelava lingvističnih podatkov*, str. 335—344, Bled, 1982.
- Guiraud, P., *Les caractères statistique du vocabulaire*, Paris, PUF, 1954.
- Parunak, H. V. D., »Prolegomena to Pictorial Concordances«, *Computers and the Humanities*, 15, str. 15—16, 1981.
- Pravda, V., »Expériences d'enseignement des langues de spécialité à l'Université Charles à Prague«, u *Les langues de spécialité analyse linguistique et recherche pédagogique*, Strasbourg, 1970.
- Teubert, W., »Corpus and Lexicography«, *Zbornik II znanstvenega srećanja računalniška obdelava lingvističnih podatkov*, str. 275—301, Bled, 1982.

⁸ Descamps, J. L. i dr.

COMPUTER CONCORDANCES IN COMPILING PEDAGOGICAL LSP DICTIONARIES

Summary

Concordances are standard tool in modern lexicology. Pedagogical LSP dictionaries are first of all functional types of dictionaries which enable the user to decode text efficiently using minimal lexical resources.

To compile the pedagogical LSP dictionary we need first of all frequency lists of the words used in a certain field (with the distribution data if possible) and concordance of the KWIC type. Using those data we have to: lemmatize the word lists, separate homographs and establish the most frequent and important collocations and idioms.

All the data should be listed by the alphabetical order, explained in the other language, chequed and presented in the final dictionary form.

Another form of the pedagogical LSP dictionary could be so called contextual dictionary of the given field, which is today almost impossible to compile without using concordances.