

Enaindvajseta mednarodna olimpijada iz jezikoslovja

Brasilia (Brazilija), 23.–31. julij 2024

Naloga skupinskega tekmovanja

Leksikostatistika je skupina metod, s katerimi lahko na podlagi besedišča ocenimo, kako tesno so si jeziki med seboj sorodni. Te metode se običajno uporabljajo skupaj z dolgimi seznammi besed, ki jih strokovnjaki ročno označijo in navedejo, ali naj bi določen par besed imel skupen izvor. Včasih pa jezikoslovci uporabljajo leksikostatistične metode za sezname besed, ki so označeni z avtomatskimi postopki. Eden takih postopkov temelji na konceptu *konzonantnih razredov*, ki ga je uvedel sovjetsko-izraelski jezikoslovec Aron Dolgopoljski leta 1964.

P.	p b ɓ φ β f v	K.	k g x γ q ɠ χ uɣ	Y.	j ɟ (na začetku korena)	M.	m ɱ
T.	t d ɗ θ ð ʈ ɖ	R.	r ɾ ɽ l ɭ ʎ ʟ ʞ	W.	w ɰ (na začetku korena)	N.	n ɲ ɳ ɰ
S.	s z ʃ ʒ ʂ ʐ ʑ ɕ ɟ					Q.	ʈ ɖ ʎ
H.	h ɣ ɦ ʕ ʔ h ɦ ʔ, samoglasniki in j ɟ w ɰ (razen na začetku korena)						

Konzonantni razredi Dolgopoljskega

Spodaj so navedeni označeni deli seznamov besed za več jezikovnih družin s celega sveta. Oznake so podane s podpisanimi številkami. Na podlagi teh seznamov so bila s pomočjo dveh poenostavljenih različic t. i. algoritma *StarlingNj* zgrajena drevesa jezikovnih družin, vsaki besedi pa je bil dodeljen *indeks stabilnosti*. Drevesa in indeksi stabilnosti na vrhu temeljijo na ročno označenih seznamih besed, tisti na dnu pa temeljijo na avtomatsko označenih seznamih. Za vsak seznam besed sta izdelani dve drevesi: eno z algoritmom A in drugo z algoritmom B. V nekaterih primerih obstaja več možnih dreves, ki ustrezajo posameznemu seznamu besed; v takih primerih je bilo naključno izbrano samo eno drevo. Vsakemu vozlišču na vsakem drevesu je pripisana leksikostatistična razdalja. Večja kot je razdalja, tesnejše je razmerje med jezikoma. Natančneje bi bilo torej uporabiti izraz „obrnjena leksikostatistična razdalja“ namesto „leksikostatistična razdalja“. Zaradi poenostavitve je v tej nalogi uporabljen izraz „leksikostatistična razdalja“.

Indeksi stabilnosti in leksikostatistične razdalje so zaokroženi na dve decimalni mesti. Če je tretja številka za decimalno vejico manjša od 5, zaokrožite navzdol, sicer zaokrožite navzgor. Npr. 2,836 se zaokroži na 2,84; 0,705 se zaokroži na 0,71; 0,703 pa se zaokroži na 0,70. Zaokroževanje velja samo za vrednosti, ki so prikazane ljudem. Z drugimi besedami, računalnik, ki poganja algoritme, „vidi“ nezaokrožene vrednosti.

Upoštevajte, da so nekatere besede gotovo ali domnevno izposojene iz drugih jezikov. Tako je beseda **jok*i*** ‘sol’ kadivejščine izposojena iz gvaranske besede **juk*i***, beseda **ʔa:n*i*** ‘leto’ ipajščine (mesa-grandsko narečje) pa je izposojena iz španske besede **ʔaño**.

V nekaterih primerih je na seznamih besed za posamezni pomen navedenih več sinonimov, ki so ločeni z vejico. Tak primer je ‘noga’ v vehoščini.

V spodnjih podatkih so vse predpone ločene z znakom „=“, vse pripone pa z znakom „-“. Nekatere besede se vedno uporabljajo s predponami. Začenjajo se z znakom „=“.

Podatki so zapisani z mednarodno fonetično abecedo. ^ˈ = glavni naglas, _ˈ = stranski naglas (šibkejši kot glavni naglas), ː = dolg glas, ˚ = zelo kratek glas, \widehat{XY} = X in Y se izgovarjata kot en glas, ˊ = visok ton, ˋ = nizek ton, ˆ = padajoč ton, ʔ˚ = predglotaliziran glas (s predhodno kratko prekinitvijo zračnega toka v grlu), ˚ʔ = ejektivni soglasnik (izgovorjen s kratko prekinitvijo zračnega toka v grlu), ˚ = nezveneč soglasnik, ˚̃ = nosni glas (izgovorjen skozi nos), ˚̥ = škripajoč glas (nizek, prasketajoč

zvok), n° označuje zračni tok skozi nos pred soglasnikom, o^h = pridihneni soglasnik (izgovori se s pridihom zraka), o^w = labializiran soglasnik (izgovorjen z zaobljenimi ustnicami), o^j = mehčan glas (izgovori se z delom jezika pri trdem nebu). $\alpha, \text{æ}, \text{ɛ}, \text{ɪ}, \text{ɨ}, \text{ɔ}, \text{ʊ}, \text{u}, \text{ə}, \text{ʌ}, \text{ɒ}, \text{ɘ}, \text{y}, \text{ø}, \text{ø}$ so samoglasniki. Drugi posebni znaki so soglasniki.

△ Znanje jezikov, omenjenih v nalogi, ne predstavlja prednosti pri reševanju naloge.

I. sklop. Gvajkurujška družina (Argentina, Brazilija, Paragvaj)

	toba (vzhodna)	pilaga	mokovijščina (čakovsko)	kadivejščina
oblak	l=ʔok ₁	'lo=ʔok ₁	naweyelek ₂	lol:adi ₃
ogenj	nodek ₁	'd=oleʔ ₂	norek ₁	n=ol:edi ₂
riba	njaq ₁	'nijaq ₁	naʎin ₂	nij:ogo-ḍʒegi ₃
glava	=qajk ₁	=ʔajk ₁	=qaik ₁	=ak:ilo ₂
ubiti	=alawat ₁	=aʔa:t ₁	=alawat ₁	=el:owadi ₁
luna	ʔawoʂojk ₁	ʔaʔwoʂojk ₁	ʃirajyo ₂	ep:enaj ₃
nos	=mik ₁	=ʔmik ₁	=mik ₁	=m:iq:o ₁
sol	towe ₁	olʔyek ₂	ʔwe ₁	jok:i ₋₁
kamen	qaʔ ₁	'qaʔ ₁	qaʔ ₁	wet:iga ₂
jezik	=atʃ-aʂat ₁	=aʔʃ-aʂat ₁	=oʔley-aʂan-aʂat ₂	=ok:el:i ₃

	algoritem A	algoritem B	
ročno			Indeks stabilnosti: oblak 0,50 ogenj 0,50 riba 0,50 glava 0,75 ubiti 1,00 luna 0,50 nos 1,00 sol 0,67 kamen 0,75 jezik 0,50
avtomatsko			Indeks stabilnosti: oblak 0,50 ogenj 0,50 riba 0,75 glava 0,75 ubiti 1,00 luna 0,50 nos 1,00 sol 0,25 kamen 0,75 jezik 0,50

II. sklop. Nubijska družina (Egipt, Sudan)

	dongolavijščina	kenuzijščina	dilingijščina	kadarujščina	debrijščina	birgidščina
ubiti	'bɛ:₁	be:₁	hur₂	wur-i₂	wur-i₂	fila:l-e₁
luna	u'n-at-t₁	an-at-ti₁	nɔn-ti₁	nɔn-tu₁	nɔn-to₁	ma:l₂
voda	'ɛss₁	essi₁	ɔti₁	ɔto₁	ɔtu₁	eji₁
dati	'tir₁	tir₁	ti₁	ti₁	ti₁	te:n₁
dober	'sɛrɛ:₁	sere:₁	ken₂	kɛn₂	kɛŋ₂	azze-n₃
veter	'turug₁	turug₁	irf-i₂	irf-o₂	irf-o₂	kurr-i₃
lasje	'dil-ti₁	si:r₂	tel-ti₁	til-tu₁	til-tu₁	ur=dill-e₁
trebuh	'tu:₁	tu:₁	te-te₂	to₁	to₁	tu:₁
spati	'nɛ:r₁	ne:r₁	jer₁	dwallɛli₂	jer-i₁	ne:r-i₁
sonce	'masil₁	masil₁	ɛj₂	aju₂	ɛŋgal-to₃	ʔi:zi₂

	algoritem A	algoritem B	
ročno			Indeks stabilnosti: ubiti 0,50 luna 0,83 voda 1,00 dati 1,00 dober 0,50 veter 0,50 lasje 0,83 trebuh 0,83 spati 0,83 sonce 0,50
avtomatsko			Indeks stabilnosti: ubiti 0,33 luna 0,50 voda 0,50 dati 0,67 dober 0,50 veter 0,50 lasje 0,83 trebuh 1,00 spati 0,50 sonce 0,50

- (A) (2 točki) Soglasnik **ɤ** se izgovarja kot francoski *r*, torej z zadnjim delom jezika. V kateri razred Dolgopoljskega spada in kako ste to ugotovili?
- (B) (2 točki) Nubijsko drevo zgoraj levo je eno od dveh možnih dreves za to kombinacijo algoritma in označevalnega načina. Narišite drugo možno drevo.
- (C) (2 točki) Nubijsko drevo spodaj levo je eno od dveh možnih dreves za to kombinacijo algoritma in označevalnega načina. Narišite drugo možno drevo.
- (D) (2 točki) Leksikostatistična razdalja 0,49, dodeljena korenu nubijskega drevesa zgoraj desno, je zaokrožena na dve decimalni mesti, tako kot nekatere druge razdalje v tej nalogi. Kakšna je točna razdalja?

III. sklop. Matagvajska družina (Argentina, Bolivija, Paragvaj)

	vičijščina (dolnjeber- meško)	vičijščina (rivada- vijsko)	vehoščina	venaješčina	ijohvaha	manhujščina	nivaklejščina (dolnjepil- komajsko)	nivaklejščina (gornjepil- komajsko)	maka
ogenj	ʔitoχ ₁	ʔitəχ ₁	ʔitah ₁	ʔi:taχ ₁	ʰwat ₂	ʔeite ₁	ʔitaχ ₁	ʔitaχ ₁	feʔt ₂
riba	ʔwahat ₁	wahat ₁	wahat ₁	ʔwa:hat ₁	siʔjus ₋₁	ʃiʔjus ₋₁	saxetʃ ₋₁	saxetʃ ₋₁	sehets ₋₁
noga	=patʃu ₁	=qəɓ ₂	=patʃo ₁ , =kala ₂	=pa:kʔo ₁	=ʔsat ₃	=kaʔla ₂	=φo ₄	=φo ₄	=fʔi ₅
voda	ʔinot ₁	ʔinət ₁	wah ₂	ʔina:t ₁	ʔiʔnat ₁	ʔaʔnat ₁	jinaʔt ₁	jinaʔt ₁	iweli ₃
dati	=ʔweŋ-u ₁	=weŋ-u ₁	=ʔweŋ-o ₁	=ʔweŋ-o ₁	=ʔwehn-a ₁ m ₂	=ʔhaj ₃ , =ʔweŋ ₂	=xut ₄	=xut-ej ₄	tis-ix ₅
dober	ʔis ₁	ʔis ₁	ʔis ₁	ʔis ₁	ʔes ₁	ʔeis ₁	ʔis ₁	ʔis ₁	t=ejkʔun-ej ₂
veter	ʔinwok ^w ₁	ʔinwək ₁	ʔihwok ^w ₁	=ja:ʔ ₂ , =x ^w ox ^w ₃	ʔhlahwu ₄	ʔhlahwu ₄	ʔaβiʔm ₅	ʔaβiʔm ₅	tʔunikʔi ₆
drevo	haʔlo ₁	halə ₁	haʔla ₁	haʔla ₁	ʔaʔla ₁	ʔaʔla-k ₁	ʔaʔkxi-juk ₂	jiʔkla ₁	naxka-k ₃
lasje	=ʔwule-j ₁	=wule-j ₁	=ʔwole-j ₁	=ʔwo:le-ç ₁ , hi:lenaχ ₂	=ʔwole ₁	=ʔwole-j ₁	=ʔateʔtʃ ₃	=jeʔs ₄	=ʔewkux-its ₅
ubiti	=lon ₁	=lən ₁	=lan ₁	=la:ŋ ₁	=ʔlaʔan ₁	=ʔlan ₁	=kla ₁	=kla ₁	=lan ₁

	algoritem A	algoritem B	
ročno			Indeksi stabilnosti: ogenj 0,78 riba 1,00 noga 0,33 voda 0,78 dati 0,44 dober 0,89 veter 0,33 drevo 0,78 lasje 0,67 ubiti 1,00
avtomatsko			Indeksi stabilnosti: ogenj 0,78 riba 0,44 noga 0,33 voda 0,56 dati 0,67 dober 0,89 veter 0,22 drevo 0,67 lasje 0,67 ubiti 1,00

IV. sklop. Mongolska družina (Ljudska republika Kitajska, Mongolija, Rusija)

(E) (10 točk) Preglejte naslednji seznam besed. Izračunajte indekse stabilnosti, ki ustrezajo ročnim in avtomatskim oznakam.

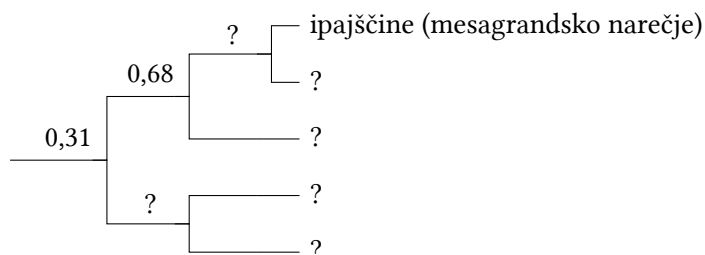
Za lažje reševanje smo za besedo 'vse' že izračunali obadva indeksa stabilnosti. V naključnem vrstnem redu sta to: 0,36 in 0,40.

	dagurščina (hajlar- sko)	hamnigan- ščina (mandžur- sko)	burjatščina (horijsko)	novo- bargut- ščina	oldščina	hošutščina	kalmiščina	halhaščina	ordoščina	širojugur- ščina	bonanščina
vse	hɔ:₁	bɔlt₂	buxi:₃	bygd₄	tsug₅	lug₅	tsuk₅, xamak₋₁	pux₃, pugt₄, xamāg₋₁	pyyite₄, xamukᵃ₋₁	tʃᵃuq₅	hanə₂
lubje	hails₁	qalihɔn₁	χoltɔhɔn₂	xalʃhu:₁	xolts₂	xalis₁	dursn₃	xɔʃtᵃs₂	turusu₃	χalsən₁	arasun₄
trebuh	ke:li₁	getəhɔn₂	gedehen₂	gedy:₂	ge:s₂	gets₂	gesn₂	gitis₂, xiwʃij₋₁	ketysy₂	ketesən₂	kele₁
ptica	dəgi₋₁	eiwan₁	ʃubu:n₁	ʃuwu:₁	ʃuvu:₁	ʃuwu:₁	ʃowun₁	ʃuwu₁	ʃuβu:₁	ʃu:n₁, peltʃər₂	bendžer₂
ogenj	gali₁	gal₁	gal₁	gal₁	gal₁	gal₁	gal₁	gal₁	qal₁	qal₁	χal₁
pot	terg-u:l₁	qargɔi₂	χargi₂, zam₋₁	zam₋₁	džam₋₁	džam₋₁	xa:-lɔə₃	tsam₋₁	tjam₋₁	mør₄	mor₄
sol	hata:₁	dawhɔn₂	dabhan₂	dawuhu:₂	daws₂	daws₂	dawsn₂	tawsā₂	taβusu₂	ta:psən₂	dabsuŋ₂
plavati	unpa-du₁	ɔmba₋₁	tᵃamar₋₂	umb₋₁	sele₋₃	umba₋₁	us-tɕi₋₄, ø:m₋₅	siʃi₋₃	usu-tʃᵃi-la₋₄	umpa₋₁	mba₋₁
voda	ɔsɔ₁	ɔxɔn₁	uhan₁	u:ha₁	usn₁	us₁	usn₁	ʊsɔ₁	usun₁	qᵃusun₁	sə₁
veter	kein₁	halkin₂	halxin₂	halxi₂	salʃxin₂	salkʃi₂	salʃkn₂	salʃxi₂	kᵃi:₁	kᵃi:₁	ki₁

V. sklop. Jumska družina (Mehika, ZDA)

(F) (8 točk) Preglejte naslednji seznam besed. Spodaj si lahko ogledate drevo, ki je bilo sestavljeno na podlagi istega seznama. Manjkajo nekateri podatki (imena jezikov in leksikostatistične razdalje). Zapolnite prazna mesta. Opredelite, ali je drevo izdelano ročno ali avtomatsko in ali je bilo generirano z algoritmom A ali B.

	mohavejščina	kokopa	javapajščina	tipajščina (hamulsko)	ipajščine (mesagrandsko narečje)
kratek	wena=wen-a ₁	'xɬ=ʔut ₂	'tʃkr=ot-i ₂	lə=ʔuj ₁	mə=put-k ₃
ptica	ʔitʃ=i=jer ₁	'ʃa ₂	'ʔ=tʃ=sa ₂	aʔ='ʃa ₂	ʔa:=ʃa:₂
kost	ɲ=a=s=ak ₁	'ɲ=j=a:k ₁	'tʃ=j=a:k-a ₁	'ak ₁	aq ₁
suh	i=ro:-v-k ₁	'ʃ=ʔar ₂	'ru-β-i ₁	's=ʔa:j ₃	sa:j ₃
meso	k ^w i:k ^w ay ₁	ʔi='ma:tʃ ₂	'k ^w e:=ʔo-β-a ₃	'k ^w ak ₄	kuk ^w a:j-p ₁
vrat	maʃaqe ₁	'm=puk ₂	'mlq ₁	i:=ʔuk ₂	i:=puk ₂
videti	i=ju:-k ₁	'wi:₂	'ʔu:₁	'wi:w ₂	ə=wu:w ₂
rep	i:=ʔar ₁	'ʃ=juʃ ₂	'β=hé ₃	ʃə='juʃ ₂	xə=juʃ ₂
dva	havik-k ₁	'x=wak ₁	'h ^w âk-i ₁	xə='wak ₁	xə=wak ₁
leto	hu:ðe ₁	'mat-'ka:m ₂	'ʔ=tʃ ^h ur-a ₃	mat-'wam ₂	ʔa:n ⁱ -₁



(G) (20 točk) Za jumske jezike so bila sestavljena še nekatera druga drevesa z naslednjimi leksikostatističnimi razdaljami v korenu drevesa (na skrajni levi strani vsakega drevesa):

1. 0,20
2. 0,23
3. 0,24

Narišite vsako od teh dreves. Za vsako drevo opredelite, ali je izdelano ročno ali avtomatsko, in ali je bilo generirano z algoritmom A ali B.

(H) (3 točke) Dve razdalji, navedeni v izzivu (G), sta zaokroženi na dve decimalni mesti: 0,23 je nastalo z zaokrožitvijo 0,225. Katera druga razdalja je bila zaokrožena in kakšna je njena točna vrednost?

(I) (4 točke) Pojasnite, kako se izračunajo indeksi stabilnosti.

(J) (5 točk) Pojasnite, kako se izračunajo leksikostatistične razdalje.

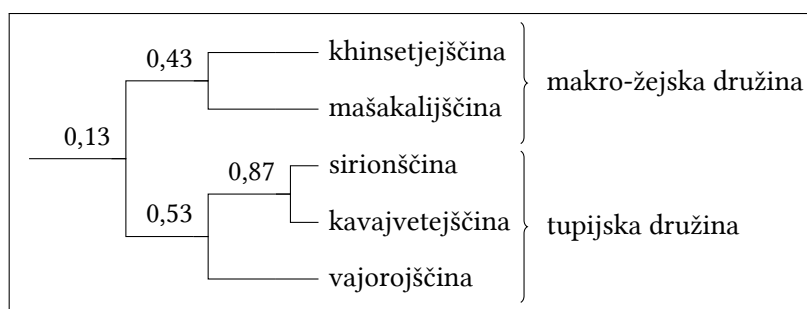
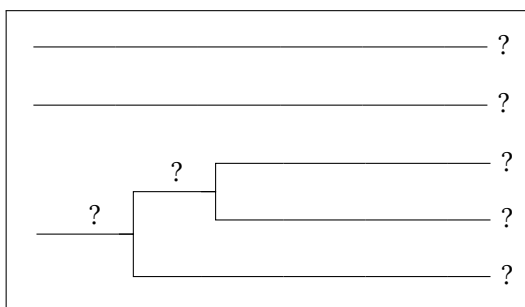
(K) (4 točke) Pojasnite razliko med algoritmom A in B.

VI. sklop. Makro-žejska družina in tupijska družina (Brazilija, Bolivija)

(L) (28 točk) Makro-žejska in tupijska družina sta pomembni jezikovni družini v Južni Ameriki. Nekateri jezikoslovci menijo, da sta v daljnem sorodstvu. Preglejte naslednje sezname besed.

	A	B	Γ	Δ	E
lubje	e='e-ke	h ^w i='k ^h Λ	kup='pε	mīβm='tεaj	= 'pε
trebuh	'e=rje	=t ^h igi	=ã'ün	= 'tæj	=rε'wεk
kri	e='ruki	=ka' ⁿ bɾo	=d̄z=a'u	= 'hεβp	=ru'i
žgati	= 'raĩ	=rɔ='k ^h ɔ̃	=po'k ^w a	mũ=...='haβp	=ra'pi
maščoba	e='kira	=t ^h wəmi	=d̄z=ap	= 'tuβp	= 'kap
noga	'e=i	= 'h ^w aji	=βi	=pɔ'ta	= 'pi
roka	'e=o	=ɲi' ^{k^h} ja	=βo	= 'ɲiβm	= 'pɔ
težek	e='usi	=wi't ^h i	=po'ti	=βp'təj	=pɔ'ij
jetra	'e=ja	= 'nba	=pi'a	=tεiβpkĩ'nāj	=pi'ʔa
nov	e='jasu	= 'ndiwi	=pa'gop	= 'tiβp	=pia'u
korenina	e='rao	=ja'rje	kup=kujopε	mīβm=ɲiβm=tεa'tiə	=ra'pɔ
koža	'e=i	= 'k ^h Λ	= 'pε	= 'tεaj	= 'pit
rep	e='rokoĩ	= 'nbi	=d̄z=o'k ^w aj	=nã:='kiβp	= 'raj
bel	'e=fi	=ja'k ^h a	=d̄zi'ra	=βp'dou	= 'sĩɲ
krilo	e='heo	=ja'rja	=pε'o	=ɲi'māu	=pε'pɔ, =ji'wa

Spodaj si lahko ogledate dve drevesi, ki sta bili sestavljeni na podlagi istih seznamov. Manjkajo nekateri podatki (imena jezikov in leksikostatistične razdalje). Zapolnite prazna mesta. Za vsako drevo opredelite, ali je izdelano ročno ali avtomatsko, in ali je bilo generirano z algoritmom A ali B.



A	B	Γ	Δ	E
?	?	?	?	?

⚠ Pri tem izzivu so bile namerno izpuščene ročno pripravljene oznake in indeksi stabilnosti.

(M) (10 točk) Avtomatizirani postopki, ki temeljijo na razredih Dolgopoljskega, lahko dajo napačne rezultate. V tem primeru avtomatski postopek odkrije več podobnosti med sirionščino in določenim makro-žejskim jezikom (khinsetjejščino) kot med sirionščino in drugimi tupijskimi jeziki. Predlagajte spremenjen avtomatiziran postopek, ki bo dal pravilno klasifikacijo, če se uporabi za zgornja makro-žejska in tupijska seznama besed, in ga *na kratko* opišite.

⚠ Ta izziv bo ocenjen le v primeru enakega števila točk med ekipami, ki so dosegle najvišje število točk.

Avtorji se zahvaljujejo Alejandri Vidal, Mariji Konošenko, Ilji Gruntovu in Jamthu Suyu za odgovore na vprašanja o določenih jezikih. —*Andrej Nikulin, Milena Veneva*

Uredniki: Ivan Deržanski (tehn. ur.), Hugh Dobbs, Stanislav Gurevič, Boris Iomdin, Eimear McKnight, Andrej Nikulin (gl. ur.), Aleksejs Peguševs, Jan Petr, Aleksandr Piperski, Marija Rubinštejn, Milena Veneva, Elisija Warner.

Slovensko besedilo: Andrej Perdih.

Srečno!