

Kahekümne esimene rahvusvaheline lingvistikaolümpiaad

Brasília (Brasília), 23.–31. juuli 2024

Meeskonnavõistluse ülesanne

Leksikostatistika on grupp meetodeid, mida kasutatakse selleks, et hinnata, kui lähedased on keeled üksteisele sõnavara poolest. Neid meetodeid kohaldatakse tavaliselt pikkadele sõnaloeteludele, mida spetsialistid käsitsi kommenteerivad, et näidata, kas antud sõnapaar on eeldatavalt samast allikast pärit. Mõnikord kasutavad keeleteadlased siiski leksikostatistilisi meetodeid automaatselt märgendatud sõnaloeteludele. Üks selline meetod põhineb *konsonandiklasside* mõistel, mille nõukogude-iisraeli keeleteadlane Aharon Dolgopolski võttis kasutusele aastal 1964.

P.	p b ɸ β f v	K.	k g x ɣ q ɕ χ ɰ	Y.	j ç (tüve alguses)	M.	m ɱ
T.	t d θ ð t̪ d̪	R.	r r̥ ɽ l ʎ ʒ ʎ ʎ	W.	w ɱ (tüve alguses)	N.	n ɲ ɳ ŋ
S.	s z ʃ ʒ ʂ ʐ ʑ ʒ ʑ					Q.	ʈ ɖ ʑ
H.	h ʕ ɦ ʕ ʔ h ɦ ʔ, vokaalid ja j ç w ɱ (välja arvatud tüve alguses)						

Dolgopolski konsonandiklassid

All leiate märkustega varustatud sõnaloetelude väljavõtted maailma erinevatest keelkondadest. Märkused on esitatud alaindeksitega. Nende loetelude põhjal on konstrueeritud keelepuud, kasutades nn *StarlingN7* algoritmi kahte lihtsustatud meetodit, ning igale sõnale on määratud *stabiilsusindeks*. Ülemises osas esitatud puud ja stabiilsusindeksid põhinevad käsitsi märgendatud sõnaloeteludel ja alumises osas esitatud puud ja stabiilsusindeksid põhinevad automaatselt märgendatud loeteludel. Iga sõnaloetelu jaoks on olemas kaks konstrueeritud puud, vastavalt algoritmi mõlemale versioonile: Algoritm A ja Algoritm B. Pange tähele, et mõnikord on mitu võimalikku puud, mis vastavad ühele sõnaloetelule; sellisel juhul on üks puu valitud juhuslikult. Iga sõlmpunktile on määratud leksikostatistiline kaugus. Mida suurem on kaugus, seda lähedasemalt on keeled omavahel suguluses. Täpsem mõiste oleks siis „vastupidine leksikostatistiline kaugus“, mitte „leksikostatistiline kaugus“. Et asja lihtsamaks teha, kasutame siin ülesandes mõistet „leksikostatistiline kaugus“.

Nii stabiilsusindeksid kui ka leksikostatistilised kaugused on ümardatud kahe kümnendkohani. Kui kolmas number pärast koma on väiksem kui 5, ümardatakse allapoole; vastasel juhul ümardatakse ülespoole. Näiteks on 2,836 ümardatud 2,84 peale, 0,705 0,71 peale ja 0,703 0,70 peale. Ümardamine kehtib ainult nende väärtuste kohta, mida näidatakse inimestest lugejatele. Seega arvuti „näeb“ algoritmide rakendamise käigus ümardamata väärtusi.

Pange tähele, et mõned sõnad on teadaolevalt või arvatavasti laenatud teistest keeltest. Näiteks kadiveu keele **jok:i** ‘sool’ sõna on laenatud guaranikeelsest **juki** ning ’iipay (Mesa Grande murre) keele **?a:nj** ‘aasta’ sõna on laenatud hispaaniakeelsest **’ano**.

Mõnel juhul on sõnade loeteludes ühe tähenduse jaoks esitatud mitu sünonüümi, mis on üksteisest eraldatud komaga. Üks näide on ‘jalg’ vehhose keeles.

Allpool esitatud andmetes on kõik eesliited eraldatud märgiga „=“ ja kõik järelliited märgiga „-“. Mõnda sõna kasutatakse ainult eesliitega koos. Need algavad märgiga „=“.

Andmed on transkribeeritud rahvusvahelise foneetilise tähestiku (IPA) abil. ^ˈ = pearõhk, _ˌ = kaasarõhk (nõrgem kui pearõhk), ː = pikk häälik, ˘ = ülilühike häälik, \widehat{XY} = X ja Y hääldatakse ühe häälikuna, ˊ = kõrge toon, ˋ = madal toon, ˆ = langev toon, ˚ = preglotaliseeritud häälik (eelneb lühike õhuvoolu katkestus kurgus), ˗ = ejetiiv (konsonant, mille hääldamisel kurk korraks sulgub), ˚ = helitu häälik, ˚ = nasaliseerunud häälik (nina kaudu hääldatav häälik), ˚ = kärisev häälik (madala, ragiseva kõlaga),

n tähendab, et enne kaashäälikut läheb osa õhust läbi nina, h = aspireeritud konsonant (hääldusega kaasneb lisahõngus), w = labialiseeritud kaashäälik (hääldatakse ümarate huultega), j = palataliseeritud (pehmendatud) häälik (hääldamisel liigub keel kõva suulae poole). **a, æ, ε, i, ɔ, u, ɥ, ə, ʌ, ɒ, ɔ, ʊ, ø, ʊ** on täishäälikud. Muud erimärgid tähistavad kaashäälikuid.

△ Mõne ülesandes kasutatud keele oskamine ei anna ülesande lahendamiseks mingeid eeliseid.

I osa. Guaikuru keelkond (Argentina, Brasiilia, Paraguay)

	toba (idamurre)	pilaga	mokovi (chaco murre)	kadiveu
pilv	l=ʔok ₁	'lo=ʔok ₁	naweyelek ₂	lol:adi ₃
tuli	nodek ₁	'd=oleʔ ₂	norek ₁	n=ol:edi ₂
kala	njaq ₁	'nijaq ₁	naʎin ₂	nij:ogo-ḍʒegi ₃
pea	=qajk ₁	'qajk ₁	=qaik ₁	=ak:ilo ₂
tapma	=alawat ₁	=a'la:t ₁	=alawat ₁	=el:owadi ₁
kuu (taevakeha)	ʔawoʂojk ₁	ʔa'woʂojk ₁	ʃirajyo ₂	ep:enaj ₃
nina	=mik ₁	'mik ₁	=mik ₁	=m:iq:o ₁
sool	towe ₁	ol'yek ₂	ʔwe ₁	jok:i ₁
kivi	qaʔ ₁	'qaʔ ₁	qaʔ ₁	wet:iga ₂
keel	=atʃ-aʂat ₁	=a'tʃ-aʂat ₁	=oʔley-aʂan-aʂat ₂	=ok:el:i ₃

	A-algoritm	B-algoritm	
käsitsi	<p>leksikostatistiline kaugus</p>		Stabiilsusindeksid: pilv 0,50 tuli 0,50 kala 0,50 pea 0,75 tapma 1,00 kuu (taevakeha) 0,50 nina 1,00 sool 0,67 kivi 0,75 keel 0,50
automatiseeritud			Stabiilsusindeksid: pilv 0,50 tuli 0,50 kala 0,75 pea 0,75 tapma 1,00 kuu (taevakeha) 0,50 nina 1,00 sool 0,25 kivi 0,75 keel 0,50

II osa. Nuubia keelkond (Egiptus, Sudaan)

	dongolawi	kenuzi	dilingi	kadaru	debri	birgidi
tapma	'bɛ:₁	be:₁	hur₂	wur-i₂	wur-i₂	fila:l-e₁
kuu (taevakeha)	u'n-at-t₁	an-at-ti₁	nɔn-ti₁	nɔn-tu₁	nɔn-to₁	ma:l₂
vesi	'ɛss₁	essi₁	ɔti₁	ɔto₁	ɔtu₁	ɛji₁
andma	'tir₁	tir₁	ti₁	ti₁	ti₁	te:n₁
hea	'sɛrɛ:₁	sere:₁	ken₂	kɛn₂	kɛŋ₂	azze-n₃
tuul	'turug₁	turug₁	irf-i₂	irf-o₂	irf-o₂	kurr-i₃
juuksed	'dil-ti₁	si:r₂	tel-ti₁	til-tu₁	til-tu₁	ur=dill-e₁
kõht	'tu:₁	tu:₁	te-te₂	to₁	to₁	tu:₁
magama	'nɛ:r₁	ne:r₁	ɟer₁	dwallɛli₂	ɟer-i₁	ne:r-i₁
päike	'masil₁	masil₁	ɛɟ₂	aju₂	ɛŋgal-to₃	ʔi:zi₂

	A-algoritm	B-algoritm	
käsitsi			tapma 0,50 kuu (taevakeha) 0,83 vesi 1,00 andma 1,00 hea 0,50 tuul 0,50 juuksed 0,83 kõht 0,83 magama 0,83 päike 0,50
automatiseeritud			tapma 0,33 kuu (taevakeha) 0,50 vesi 0,50 andma 0,67 hea 0,50 tuul 0,50 juuksed 0,83 kõht 1,00 magama 0,50 päike 0,50

- (A) (2 punkti) Kaashäälik **ɟ** hääldub nagu prantsuse *r* ehk keele tagaosas. Millisesse Dolgopolski klassi see kuulub ja kuidas te teada saite?
- (B) (2 punkti) Nuubia keelte puu üleval vasakul on ainult üks kahest võimalikust puust selle algoritmi ja märkustetüübi kombinatsiooni puhul. Joonistage teine võimalik puu.
- (C) (2 punkti) Nuubia keelte puu all vasakul on ainult üks kahest võimalikust puust selle algoritmi ja märkustetüübi kombinatsiooni puhul. Joonistage teine võimalik puu.
- (D) (2 punkti) Leksikostatistiline kaugus 0,49 (mis on määratud Nubia puu juurele üleval paremal) on ümardatud kahe kümnendkohani, nagu ka mõned teised selle ülesande kaugused. Mis on täpne kaugus?

III osa. Mataguai keelkond (Argentina, Boliivia, Paraguay)

	vitši (Bermejo alamjooksu murre)	vitši (Rivadavia murre)	vehhose	'weenhayeki	iyojwa'aja'	manjui	nivakle (shichaam lhavos)	nivakle (chisham-nee lhavos)	maka
tuli	ʔitox ₁	ʔitox ₁	ʔitah ₁	ʔi:tax ₁	'hwat ₂	'ʔeite ₁	ʔitax ₁	ʔitax ₁	feʔt ₂
kala	ʔwahat ₁	wahat ₁	wahat ₁	ʔwa:hat ₁	si'ʔjus ₋₁	ʃi'ʔjus ₋₁	saxetʃ ₋₁	saxetʃ ₋₁	sehets ₋₁
jalg	=patʃu ₁	=qol ₂	=patʃo ₁ , =kala ₂	=pa:k'oʔ ₁	=sat ₃	=ka'laʔ ₂	=φoʔ ₄	=φoʔ ₄	=f'iʔ ₅
vesi	ʔinot ₁	ʔinot ₁	wah ₂	ʔina:t ₁	ʔi'n'at ₁	ʔa'ʔnat ₁	jinaʔt ₁	jinaʔt ₁	iweliʔ ₃
andma	=ʔweŋ-u ₁	=weŋ-u ₁	=ʔweŋ-o ₁	=ʔweŋ-oʔ ₁	=wɛhn-aʔm ₂	=hajʔ ₃ , =wɛn ₂	=xut ₄	=xut-ej ₄	tis-ix ₅
hea	ʔis ₁	ʔis ₁	ʔis ₁	ʔis ₁	'ʔes ₁	'ʔeis ₁	ʔis ₁	ʔis ₁	t=ejk'un-ej ₂
tuul	ʔinwok ^w ₁	ʔinwək ₁	ʔihwok ^w ₁	=ja:ʔ ₂ , =x ^w ox ^w ₃	'hlahwuʔ ₄	'hlahwu ^w ʔ ₄	ʔaβiʔm ₅	ʔaβiʔm ₅	t'unik'i ₆
puu	haʔlo ₁	hal ₁	haʔla ₁	haʔlaʔ ₁	ʔa'ʔlaʔ ₁	ʔa'ʔla-k ₁	ʔaʔkxi-juk ₂	jiʔklaʔ ₁	naxka-k ₃
juuksed	=ʔwule-j ₁	=wule-j ₁	=ʔwole-j ₁	=ʔwo:le-ç ₁ , hi:lənax ₂	=ʔwole ₁	=ʔwole-j ₁	=ʃateʔtj ₃	=jeʔs ₄	=ʔewkux-its ₅
tapma	=lon ₁	=lən ₁	=lan ₁	=la:ŋ ₁	=laʔan ₁	=lan ₁	=klan ₁	=klan ₁	=lan ₁

	A-algoritm	B-algoritm	
käsitsi			tuli 0,78 kala 1,00 jalg 0,33 vesi 0,78 andma 0,44 hea 0,89 tuul 0,33 puu 0,78 juuksed 0,67 tapma 1,00
automatiseeritud			tuli 0,78 kala 0,44 jalg 0,33 vesi 0,56 andma 0,67 hea 0,89 tuul 0,22 puu 0,67 juuksed 0,67 tapma 1,00

IV osa. Mongoli keelkond (Hiina Rahvavabariik, Mongoolia, Venemaa)

(E) (10 punkti) Vaadake järgnevat sõnaloetelu. Arvutage stabiilsusindeks, nii käsitsi tehtud kui ka automatiseeritud märkuste jaoks.

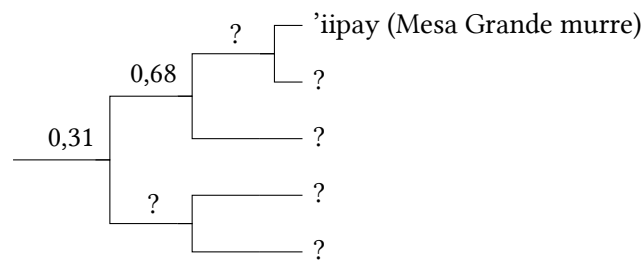
Et teid pisut aidata, oleme sõna 'kõik' jaoks mõlemad stabiilsusindeksid juba välja arvutanud. Need on – suvalises järjekorras – 0,36 ja 0,40.

	daguuri (hailari muure)	hamnigani (mandžuu- ria murre)	burjaadi (hori murre)	uusbarga	ööldi	hošuudi	kalmõki	halha	ordose	šira- juguri	bonani
kõik	hɔ:₁	bʊlt₂	bʊxi:₃	bygd₄	tsug₅	lug₅	tsuk₅, xamak₋₁	pux₃, pugt₄, xamǎg₋₁	pyyyte₄, xamukᵃ₋₁	tʃᵃuq₅	hanə₂
puukoor	hails₁	qalihon₁	χoltɔhɔn₂	xalʃhu:₁	xolts₂	xalis₁	dursn₃	xɔɣtᵃɔs₂	turusu₃	χalsən₁	arasun₄
kõht	ke:li₁	gətəhən₂	gedehen₂	gedy:₂	ge:s₂	gets₂	gesn₂	gitis₂, xiwɣij₋₁	ketysy₂	ketesən₂	kele₁
lind	dəgi₋₁	eiwan₁	ʃubu:n₁	ʃuwu:₁	ʃuvu:₁	ʃuwu:₁	ʃowun₁	ʃuwu₁	ʃuβu:₁	ʃu:n₁, peltʃɔr₂	bendzer₂
tuli	gali₁	gal₁	gal₁	gal₁	gal₁	gal₁	gal₁	gal₁	qal₁	qal₁	χal₁
rada	terg-u:l₁	qargʊi₂	χargi₂, zam₋₁	zam₋₁	dzam₋₁	dzam₋₁	xa:-lɔə₃	tsam₋₁	tʃam₋₁	mør₄	mor₄
sool	hata:₁	dawhɔn₂	dabhan₂	dawuhu:₂	daws₂	daws₂	dawsn₂	tawsǎ₂	taβusu₂	ta:psən₂	dabsuŋ₂
ujuma	unpa-du₁	ɔmba₋₁	tᵃamar₋₂	umb₋₁	sele₋₃	umba₋₁	us-tɛi₋₄, ø:m₋₅	siɣi₋₃	usu-tʃᵃi-la₋₄	umpa₋₁	mba₋₁
vesi	ɔsɔ₁	ɔxɔn₁	uhan₁	u:ha₁	usn₁	us₁	usn₁	ʊšö₁	usun₁	qᵃusun₁	sə₁
tuul	kein₁	halkin₂	halxin₂	halxi₂	salʃxin₂	salkʃi₂	salʃkn₂	saɣxi₂	kᵃi:₁	kᵃi:₁	ki₁

V osa. Juma keelkond (Mehhiko, Ameerika Ühendriigid)

(F) (8 punkti) Vaadake järgnevat sõnaloetelu. All näete puud, mis konstrueeriti sama loetelu alusel. Mõningad andmed (keelte nimed ja leksikostatistilised kaugused) on puudu. Täitke lüngad. Pange kirja, kas puu on käsitsi tehtud või automaatne, samuti, kas see on tehtud algoritmi A või B järgi.

	mohave	kokopa	javapai	tiipai (jamuli)	'iipay (Mesa Grande murre)
lühike	wena=wen-a ₁	'xʌ=ʔut ₂	'tʃkr=ot-i ₂	lə=ʔuŋ ₁	mə=put-k ₃
lind	ʔitʃ=i=jer ₁	'ʃa ₂	'ʔ=tʃ=sa ₂	aʔ='ʃa ₂	ʔa:=ʃa:2
luu	ŋ=a=s=ak ₁	'ŋ=j=a:k ₁	'tʃ=j=a:k-a ₁	'ak ₁	aq ₁
kuiiv	i=ro:-v-k ₁	'ʃ=ʔar ₂	'ru-β-i ₁	's=ʔa:j ₃	sa:j ₃
liha	k ^w i:k ^w ay ₁	ʔi='ma:tʃ ₂	'k ^w e:='θo-β-a ₃	'k ^w ak ₄	kuk ^w a:j-p ₁
kael	maʎaqe ₁	'm=puk ₂	'mlq ₁	i:='puk ₂	i:=puk ₂
nägema	i=ju:-k ₁	'wi:2	'ʔu:1	'wi:w ₂	ə=wu:w ₂
saba	i:=ʔar ₁	'ʃ=juʎ ₂	'β=hé ₃	ʃə='juʎ ₂	xə=juʎ ₂
kaks	havik-k ₁	'x=wak ₁	'h ^w âk-i ₁	xə='wak ₁	xə=wak ₁
aasta	hu:ðe ₁	'mat-'ka:m ₂	'ʔ=tʃ ⁿ ur-a ₃	mat-'wam ₂	ʔa:n ^j .1



(G) (20 punkti) Juma keelte jaoks on loodud veel mõned puud, kus puude juurte vahel on järgmised leksikostatistilised kaugused (iga puu vasakpoolses otsas):

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

Joonistage neist igaühe kohta puu. Pange iga puu kohta kirja, kas see on käsitsi tehtud või automaatne, samuti, kas see on tehtud algoritmi A või B järgi.

(H) (3 punkti) Kaks osaülesandes (G) esitatud kaugust on ümardatud kahe kümnendkohani: 0,23 on ümardatud 0,225 pealt. Milline teine kaugus on ümardatud ning mis on selle täpne väärtus?

(I) (4 punkti) Selgitage, kuidas arvutatakse stabiilsusindeksid.

(J) (5 punkti) Selgitage, kuidas arvutatakse leksikostatistilised kaugused.

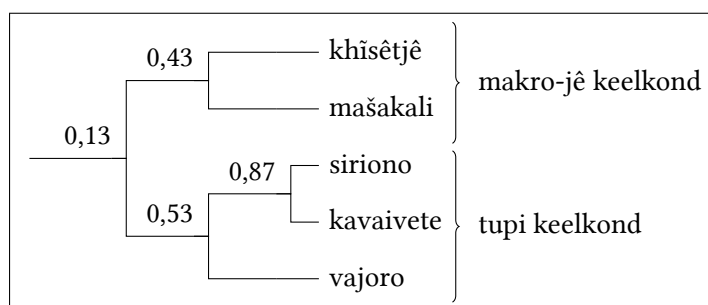
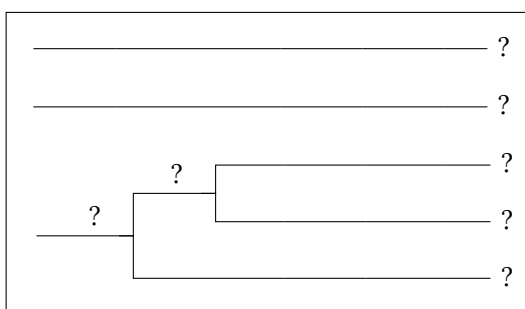
(K) (4 punkti) Selgitage A- ja B-algoritmide vahelist erinevust.

VI osa. Makro-jê keelkond ja tupi keelkond (Brasiilia, Boliivia)

(L) (28 punkti) Makro-jê ja tupi on kaks suuremat keelkonda Lõuna-Ameerikas. Mõned keeleteadlased usuvad, et nad on omavahel kaugemalt suguluses. Vaadake järgnevaid sõnaloetelusi.

	A	B	Γ	Δ	E
puukoor	e='e-ke	h ^w ĩ='k ^h λ	kup='pe	mĩβ̃m='tɛaj	'pe
kõht	'e=rje	't ^h igi	=ã'ün	'tæj	=rɛ'wek
veri	e='ruki	=ka ⁿ bɔ	=d̃z=a'u	'hɛβ̃p	=ru'i
põletama	'raĩ	=rɔ='k ^h λɔ	=po'k ^w a	mũ=...'haβ̃p	=ra'pi
rasv	e='kira	't ^h wəmi	'd̃z=ap	'tuβ̃p	'kap
jalg	'e=i	'h ^w aji	'βi	=pɔ'ta	'pi
käsi	'e=o	=ɲĩ'k ^h λa	'βo	'ɲĩβ̃m	'pɔ
raske	e='usi	=wi't ^h ĩ	=po'ti	=β̃p'tɛj	=pɔ'ij
maks	'e=ja	' ⁿ ba	=pi'a	=tɛiβ̃pkĩ'nãj	=pi'ʔa
uus	e='jasu	' ⁿ diwi	=pa'gop	'tiβ̃p	=pia'u
juur	e='rao	=ja'ɕe	kup=kujɔ'pe	mĩβ̃m=ɲĩβ̃m='tɛa'tiə	=ra'pɔ
nahk	'e=i	'k ^h λ	'pe	'tɛaj	'pit
saba	e='rokoĩ	' ⁿ bi	=d̃z=o'k ^w aj	=nã:='kiβ̃p	'raj
valge	'e=fi	=ja'k ^h a	=d̃zi'ra	=β̃p'douɥ	'sĩɲ
tiib	e='heo	=ja'ɕa	=pe'o	=ɲĩ'mãuɥ	=pe'pɔ, =ji'wa

All näete kahte puud, mis konstrueeriti samade loetelude alusel. Mõningad andmed (keelte nimed ja leksikostatistilised kaugused) on puudu. Täitke lüngad. Pange iga puu kohta kirja, kas see on käsitsi tehtud või automaatne, samuti, kas see on tehtud algoritmi A või B järgi.



A	B	Γ	Δ	E
?	?	?	?	?

⚠ Käsitsi tehtud märkused ning stabiilsusindeksid on siin osaülesandes meelega ära jäetud.

(M) (10 punkti) Dolgopolski klassidel põhinevad automaatsed protsessid võivad viia valede tulemusteni. Selles näites tuvastab automaatne protsess rohkem sarnasusi siriono ja ühe teatud makro-jê keele (khîsêjtê keele) vahel kui siriono ja teiste tupi keelte vahel. Tehke ettepanek modifitseeritud automaatse protsessi kohta, mis annaks makro-jê ja tupi sõnastike puhul õige klassifikatsiooni, ja kirjeldage seda *lühidalt*.

⚠ Seda osaülesannet hinnatakse ainult sel juhul, kui parimate tulemustega võistkonnad lõpetavad viigiga.

Autorid tänavad Alejandro Vidali, Maria Konošenkot, Ilja Gruntovit ja Jamthô Suyat, kes lahkelt vastasid küsimustele üksikute keelte kohta. —Andrei Nikulin, Milena Veneva

Toimetajad: Ivan Deržanski (tehn. toim.), Hugh Dobbs, Stanislav Gurevitš, Boriss Iomdin, Liam McKnight, Andrei Nikulin (vast. toim.), Aleksejs Peguševs, Jan Petr, Aleksandr Piperski, Marija Rubinštein, Milena Veneva, Elysia Warner.

Eesti tekst: Axel Jagau.

Edu!