

Vingt-et-unièmes Olympiades internationales de linguistique

Brasília (Brésil), du 23 au 31 juillet 2024

Problème de la compétition en équipe

La lexicostatistique est un groupe de méthodes créées pour estimer la proximité entre différentes langues à base de leurs vocabulaires. Ces méthodes sont normalement appliquées à de longues listes de mots annotées par des experts, qui indiquent si l'on pense qu'une paire de mots spécifique provient de la même source. Cependant, les linguistes appliquent parfois des méthodes lexicostatistiques aux listes de mots annotées au moyen de processus automatisés. Un tel processus est basé sur le concept des *classes de consonnes*, introduit en 1964 par le linguiste soviéto-israélien Aharon Dolgopolsky.

P.	p b ɓ φ β f v	K.	k g x ɣ q ɠ χ ɰ	Y.	j ç (au début de la racine)	M.	m ɱ
T.	t d ɗ θ ð ʈ ɖ	R.	r ɾ ɽ ɺ ɻ ʂ ʃ ʄ	W.	w ɱ (au début de la racine)	N.	n ɲ ɳ ŋ
S.	s z ʃ ʒ ʂ ʃ ʄ ʅ					Q.	ʈ ɖ
H.	h ʕ ɦ ʁ ʔ h ɦ ʔ, les voyelles, et j ç w ɱ (sauf au début de la racine)						

Les classes de consonnes de Dolgopolsky

Ci-dessous se trouvent des fragments annotés de listes de mots de quelques familles de langues du monde. Les annotations sont données au moyen de chiffres en indice. Sur la base de ces listes, des arbres généalogiques des langues ont été construits en utilisant deux versions simplifiées de l'algorithme appelé *StarlingNj*, et un *indice de stabilité* a été assigné à chaque mot. Les arbres et indices de stabilité en haut sont basés sur des listes de mots annotées manuellement, et ceux en bas sont basés sur des listes qui ont été annotées automatiquement. Deux arbres ont été construits pour chaque liste de mots, en utilisant deux versions de l'algorithme : l'Algorithme A et l'Algorithme B. Faites attention que dans certains cas il y a plusieurs arbres possibles qui correspondent à une seule liste de mots; dans ces cas un seul arbre a été choisi aléatoirement. Une distance lexicostatistique a été assignée à chaque nœud de chaque arbre. Plus la distance est grande, plus la proximité entre les langues est immédiate. Un terme plus précis est donc «distance lexicostatistique inverse» plutôt que «distance lexicostatistique». Pour simplifier, le terme «distance lexicostatistique» est utilisé dans ce problème.

Les indices de stabilité ainsi que les distances lexicostatistiques ont été arrondis à deux décimales. Si le troisième chiffre après la virgule est inférieur à 5, arrondissez à la valeur inférieure; sinon, arrondissez à la valeur supérieure. Par exemple, 2,836 est arrondi à 2,84, 0,705 à 0,71, et 0,703 à 0,70. L'arrondi ne s'applique qu'aux valeurs présentées aux lecteurs humains. Autrement dit l'ordinateur qui exécute les algorithmes peut «voir» les valeurs non-arrondies.

Notez que l'on sait ou estime que quelques mots ont été empruntés aux autres langues. Par exemple, le mot **jok:i** 'sel' de la langue kadiwéu est emprunté au mot guarani **juki**, et **?a:nj** 'année' du ipai (de Mesa Grande) est emprunté au mot espagnol **'ajno**.

Dans certains cas plusieurs synonymes d'un seul sens sont donnés dans les listes de mots, séparés par des virgules. Un exemple est 'pied' de la langue vejoj.

Dans les données ci-dessous, tous les préfixes sont séparés par le signe «=» et tous les suffixes sont séparés par le signe «-». Certains mots ne se trouvent qu'avec des préfixes. Ces mots commencent par le signe «=».

Les données ont été transcrites dans l'Alphabet Phonétique International. ' = accent primaire, , = accent secondaire (plus faible que l'accent primaire), ◌: = son long, ◌◌ = son très court, X̄Y = X et Y sont prononcés comme un seul son, ◌́ = ton haut, ◌̀ = ton bas, ◌̂ = ton descendant, ◌̚ = son préglottalisé

(précédé par un bref blocage du passage d'air dans la gorge), $\text{○}'$ = son éjectif (prononcé en bloquant brièvement le passage d'air dans la gorge), ○ = son sourd, $\text{○}̃$ = son nasalisé (prononcé avec passage de l'air par le nez), ○ = voix craquée (tonalité basse et rauque), ᵐ○ indique le flux d'air par le nez avant la consonne, ○^h = consonne aspirée (prononcée avec une bouffée d'air), ○^w = consonne labialisée (prononcée avec lèvres arrondies), ○^j = son palatalisé (prononcé en dirigeant une partie de la langue près du palais dur). $\text{a, æ, ε, i, ɪ, ɔ, ʊ, u, ə, ʌ, ɒ, ɐ, y, ø}$ sont des voyelles. Les autres caractères spéciaux sont des consonnes.

⚠ Toute connaissance des langues mentionnées dans le problème ne donne aucun avantage pour résoudre le problème.

I^{ère} partie. La famille waykuruane (Argentine, Brésil, Paraguay)

	toba (de l'Est)	pilagá	mocoví (de Chaco)	kadiwéu
nuage	l=ʔok ₁	'lo=ʔok ₁	naweyelek ₂	lol:adi ₃
feu	nodek ₁	'd=oleʔ ₂	norek ₁	n=ol:edi ₂
poisson	njaq ₁	'nijaq ₁	naʎin ₂	nij:ogo-dʒegi ₃
tête	=qajk ₁	=qajk ₁	=qaik ₁	=ak:ilo ₂
tuer	=alawat ₁	=a'la:t ₁	=alawat ₁	=el:owadi ₁
lune	ʔawoʝok ₁	ʔa'woʝok ₁	ʃirajyo ₂	ep:enaj ₃
nez	=mik ₁	=mik ₁	=mik ₁	=m:iq:o ₁
sel	towe ₁	ol'ʒek ₂	ʔwe ₁	jok:i ₁
Pierre	qaʔ ₁	'qaʔ ₁	qaʔ ₁	wet:iga ₂
langue	=atʃ-aʎat ₁	=a'tʃ-aʎat ₁	=oʔley-aʎan-aʎat ₂	=ok:el:i ₃

	algorithme A	algorithme B	
manuel	<p>distance lexicostatistique</p>		Indices de stabilité: nuage 0,50 feu 0,50 poisson 0,50 tête 0,75 tuer 1,00 lune 0,50 nez 1,00 sel 0,67 pierre 0,75 langue 0,50
automatique			Indices de stabilité: nuage 0,50 feu 0,50 poisson 0,75 tête 0,75 tuer 1,00 lune 0,50 nez 1,00 sel 0,25 pierre 0,75 langue 0,50

II^{ème} partie. La famille nubienne (Egypte, Soudan)

	dongolawi	kenzi	dilling	kadaru	debri	birgid
tuer	'bɛ:₁	be:₁	hur₂	wur-i₂	wur-i₂	fila:l-e₁
lune	u'n-at-t₁	an-at-ti₁	nɔn-ti₁	nɔn-tu₁	nɔn-to₁	ma:l₂
eau	'ɛss₁	essi₁	ɔti₁	ɔto₁	ɔtu₁	ɛj₁
donner	'tir₁	tir₁	ti₁	ti₁	ti₁	te:r-n₁
bon	'sɛrɛ:₁	sere:₁	ken₂	kɛn₂	kɛɲ₂	azze-n₃
vent	'turug₁	turug₁	irf-i₂	irf-o₂	irf-o₂	kurr-i₃
cheveux	'dil-ti₁	si:r₂	tel-ti₁	til-tu₁	til-tu₁	ur=dill-e₁
ventre	'tu:₁	tu:₁	te-te₂	to₁	to₁	tu:₁
dormir	'nɛ:r₁	ne:r₁	jer₁	dwallɛli₂	jer-i₁	ne:r-i₁
soleil	'masil₁	masil₁	ɛj₂	aju₂	ɛɲgal-to₃	ʔi:zi₂

	algorithme A	algorithme B	
manuel			Indices de stabilité : tuer 0,50 lune 0,83 eau 1,00 donner 1,00 bon 0,50 vent 0,50 cheveux 0,83 ventre 0,83 dormir 0,83 soleil 0,50
automatique			Indices de stabilité : tuer 0,33 lune 0,50 eau 0,50 donner 0,67 bon 0,50 vent 0,50 cheveux 0,83 ventre 1,00 dormir 0,50 soleil 0,50

- (A) (2 points) La consonne **ɣ** est prononcée comme l'*r* grasseyé du français, c'est-à-dire à l'arrière de la langue. À quelle classe Dolgopolsky appartient-elle, et comment l'avez-vous déterminé ?
- (B) (2 points) L'arbre nubien en haut à gauche est juste un des deux arbres possibles pour cette combinaison d'algorithme et de type d'annotation. Dessinez l'autre arbre possible.
- (C) (2 points) L'arbre nubien en bas à gauche est juste un des deux arbres possibles pour cette combinaison d'algorithme et de type d'annotation. Dessinez l'autre arbre possible.
- (D) (2 points) Comme quelques autres distances dans ce problème, la distance lexicostatistique 0,49 (assignée à la racine de l'arbre nubien en haut à droite) a été arrondie à deux décimales. Quelle est la distance exacte ?

III^{ème} partie. La famille mataguayo (Argentine, Bolivie, Paraguay)

	wichi (du Bas Bermejo)	wichi (de Rivadavia)	vejoz	'weenhayek	iyojwa'aja'	manhui	nivaklé (des Shichaam Lhavos)	nivaklé (des Chishamnee Lhavos)	maká
feu	ʔitox ₁	ʔitox ₁	ʔitah ₁	ʔi:tax ₁	ʔh ^h wat ₂	ʔe ^h it'e ₁	ʔitax ₁	ʔitax ₁	feʔt ₂
poisson	ʔwahat ₁	wahat ₁	wahat ₁	ʔwa:hat ₁	siʔʔjus ₋₁	ʃiʔʔjus ₋₁	saxetʃ ₋₁	saxetʃ ₋₁	sehets ₋₁
ped	=patʃ ^h u ₁	=qolɔ ₂	=patʃ ^h o ₁ , =kala ₂	=pa:kʔ'oʔ ₁	=ʔsat ₃	=kaʔʔlaʔ ₂	=ʃoʔ ₄	=ʃoʔ ₄	=fʔiʔ ₅
eau	ʔinot ₁	ʔinot ₁	wah ₂	ʔina:t ₁	ʔi'nʔat ₁	ʔaʔʔnat ₁	jinaʔt ₁	jinaʔt ₁	iweliʔ ₃
donner	=ʔweŋ _o -u ₁	=weŋ _o -u ₁	=ʔweŋ _o -o ₁	=ʔweŋ _o -oʔ ₁	=ʔweh ^h n-aʔm ₂	=ʔhajʔ ₃ , =ʔweŋ ₂	=xut ₄	=xut-ej ₄	tis-ix ₅
bon	ʔis ₁	ʔis ₁	ʔis ₁	ʔis ₁	ʔes ₁	ʔe ^h is ₁	ʔis ₁	ʔis ₁	t=ejk'un-ej ₂
vent	ʔinwok ^w ₁	ʔinwək ₁	ʔihwok ^w ₁	=ja:ʔ ₂ , =x ^w ox ^w ₃	ʔhlah ^h wuʔ ₄	ʔhlah ^h wuuʔ ₄	ʔaβiʔm ₅	ʔaβiʔm ₅	t'unik'i ₆
arbre	haʔlo ₁	halɔ ₁	haʔla ₁	haʔlaʔ ₁	ʔaʔʔlaʔ ₁	ʔaʔʔla-k ₁	ʔaʔkxi-juk ₂	jiʔklaʔ ₁	naxka-k ₃
cheveux	=ʔwule-j ₁	=wule-j ₁	=ʔwole-j ₁	=ʔwo:le-ç ₁ , hi:lenaç ₂	=ʔwole ₁	=ʔwole-j ₁	=ʔateʔʔʃ ₃	=jeʔs ₄	=ʔewkux-its ₅
tuer	=lon ₁	=lɔn ₁	=lan ₁	=la:ŋ ₁	=ʔlaʔan ₁	=ʔlan ₁	=kla ^h an ₁	=kla ^h an ₁	=lan ₁

	algorithme A	algorithme B	
manuel	<p> wichi (du Bas Bermejo) wichi (de Rivadavia) vejoz 'weenhayek iyojwa'aja' manhui nivaklé (des Shichaam Lhavos) nivaklé (des Chishamnee Lhavos) maká </p>	<p> wichi (du Bas Bermejo) wichi (de Rivadavia) vejoz 'weenhayek iyojwa'aja' manhui nivaklé (des Shichaam Lhavos) nivaklé (des Chishamnee Lhavos) maká </p>	Indices de stabilité : feu 0,78 poisson 1,00 pied 0,33 eau 0,78 donner 0,44 bon 0,89 vent 0,33 arbre 0,78 cheveux 0,67 tuer 1,00
automatique	<p> wichi (du Bas Bermejo) wichi (de Rivadavia) vejoz 'weenhayek iyojwa'aja' manhui nivaklé (des Shichaam Lhavos) nivaklé (des Chishamnee Lhavos) maká </p>	<p> wichi (du Bas Bermejo) wichi (de Rivadavia) vejoz 'weenhayek iyojwa'aja' manhui nivaklé (des Shichaam Lhavos) nivaklé (des Chishamnee Lhavos) maká </p>	Indices de stabilité : feu 0,78 poisson 0,44 pied 0,33 eau 0,56 donner 0,67 bon 0,89 vent 0,22 arbre 0,67 cheveux 0,67 tuer 1,00

IV^{ème} partie. La famille mongolique (République populaire de Chine, Mongolie, Russie)

(E) (10 points) Examinez la liste de mots suivante. Calculez les indices de stabilité qui correspondent aux annotations manuelles et automatisées.

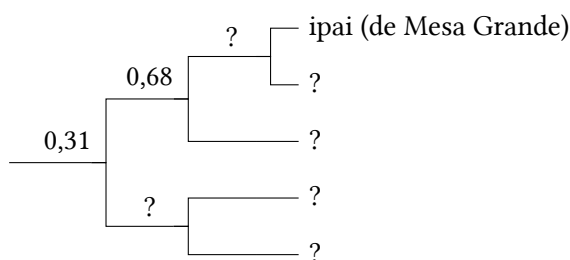
Pour vous aider, les deux indices de stabilité du mot 'tous' ont déjà été calculés. Dans un ordre aléatoire ils sont : 0,36 et 0,40.

	daur (de Hailar)	khamnigan (mand-chou)	bouriate (de Khori)	nouveau-bargu	ööld	khochoute	kalmouk	khalkha	ordos	shira yugur	bonan
tous	hɔ:₁	bölt₂	buxi:₃	bygd₄	ṭsug₅	lug₅	ṭsuk₅, xamak₋₁	pux₃, pugt₄, xamăg₋₁	pyyyte₄, xamukʰ₋₁	ṭʰuq₅	hanə₋₂
écorce	hails₁	qalihön₁	χoltöhön₂	xalʰhu:₁	xolts₂	xalis₁	dursn₃	xəɮtʰɔs₂	turusu₃	χalsən₁	arasun₄
ventre	ke:li₁	getəhən₂	gedehen₂	gedy:₂	ge:s₂	gets₂	gesn₂	gitis₂, xiwɮjij₋₁	ketysy₂	ketesən₂	kele₁
oiseau	dəgi₋₁	eiwan₁	ʃubu:n₁	ʃuwu:₁	ʃuvu:₁	ʃuwu:₁	ʃowun₁	ʃuwu₁	ʃuβu:₁	ʃu:n₁, peltʃər₂	bendzər₂
feu	gali₁	gal₁	gal₁	gal₁	gal₁	gal₁	gal₁	gal₁	qal₁	qal₁	χal₁
chemin	terg-u:l₁	qargūi₂	χargi₂, zam₋₁	zam₋₁	ḏzam₋₁	ḏzam₋₁	xa:-lɤə₃	ṭsam₋₁	ṭjam₋₁	mør₄	mor₄
sel	hata:₁	dawhön₂	dabhan₂	dawuhu:₂	daws₂	daws₂	dawsn₂	tawsă₂	taβusu₂	ta:psən₂	dabsuŋ₂
nager	unpa-du₁	umba₋₁	tʰamar₋₂	umb₋₁	sele₋₃	umba₋₁	us-təi₋₄, ø:m₋₅	siɮi₋₃	usu-tʰi-la₋₄	umpa₋₁	mba₋₁
eau	ɔsɔ₁	uxön₁	uhan₁	u:ha₁	usn₁	us₁	usn₁	ʊsö₁	usun₁	qʰusun₁	sə₁
vent	kein₁	halkin₂	halxin₂	halxi₂	salxin₂	salkʰi₂	salʰkn₂	saɮxi₂	kʰi:₁	kʰi:₁	ki₁

V^{ème} partie. La famille yumane (Mexique, États-Unis)

(F) (8 points) Examinez la liste de mots suivante. Ci-dessous est un arbre construit en utilisant la même liste. Certaines données (noms de langues et distances lexicostatistiques) manquent. Remplissez les trous. Précisez si l'arbre est manuel ou automatisé, ainsi que s'il a été généré en utilisant l'Algorithme A ou B.

	mojave	cocopa	yavapai	tipai (de Jamul)	ipai (de Mesa Grande)
court	wena=wen-a ₁	'xɬ=ʔut ₂	'tʃkr=ot-i ₂	lə=ʔuj ₁	mə=put-k ₃
oiseau	ʔitʃ=i=jer ₁	'ʃa ₂	'ʔ=tʃ=sa ₂	aʔ=ʃa ₂	ʔa:=ʃa:₂
os	ʃn=a=s=ak ₁	'ʃn=j=a:k ₁	'tʃ=j=a:k-a ₁	'ak ₁	aq ₁
sec	i=ro:-v-k ₁	'ʃ=ʔar ₂	'ru-β-i ₁	's=ʔa:j ₃	sa:j ₃
viande	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
cou	maʃaqe ₁	'm=puk ₂	'mlq ₁	i:=ʔuk ₂	i:=puk ₂
voir	i=ju:-k ₁	'wi:₂	'ʔu:₁	'wi:w ₂	ə=wu:w ₂
queue	i:=ʔar ₁	'ʃ=juʃ ₂	'β=hé ₃	ʃə=ʔuʃ ₂	xə=juʃ ₂
deux	havik-k ₁	'x=wak ₁	'h ^w âk-i ₁	xə='wak ₁	xə=wak ₁
année	hu:ðe ₁	'mat-'ka:m ₂	'ʔ=tʃ ^h ur-a ₃	mat-'wam ₂	ʔa:n ^l -₁



(G) (20 points) Quelques autres arbres ont été générés pour la famille yumane, avec les distances lexicostatistiques suivantes à la racine de l'arbre (les distances lexicostatistiques tout à gauche de chaque arbre):

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

Dessinez chacun de ces arbres. Pour chacun des arbres, précisez s'il est manuel ou automatisé, ainsi que s'il a été généré en utilisant l'Algorithme A ou B.

(H) (3 points) Deux distances données dans la tâche (G) ont été arrondies à deux décimales : 0,23 a été arrondi à partir de 0,225. Quelle autre distance a été arrondie, et quelle est sa valeur exacte ?

(I) (4 points) Expliquez comment sont calculés les indices de stabilité.

(J) (5 points) Expliquez comment sont calculées les distances lexicostatistiques.

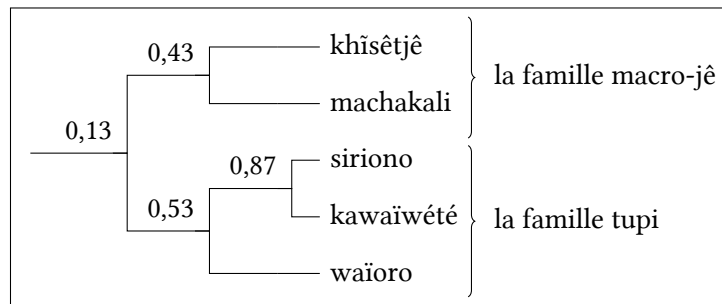
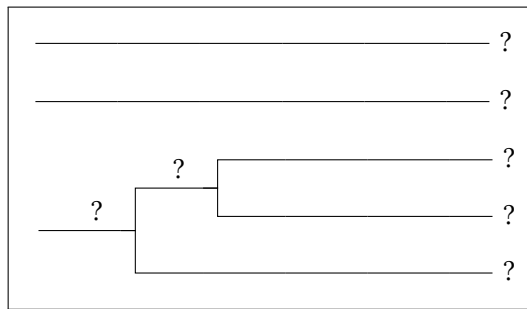
(K) (4 points) Expliquez la différence entre les algorithmes A et B.

VI^{ème} partie. La famille macro-jê et la famille tupi (Brésil, Bolivie)

(L) (28 points) La famille macro-jê et la famille tupi sont deux familles de langues importantes d'Amérique du Sud. Quelques linguistes croient qu'elles sont lointainement apparentées. Examinez les listes de mots suivantes.

	A	B	Γ	Δ	E
écorce	e='e-ke	h ^w i='k ^h Λ	kup='pe	mĩβm='tɛaj	= 'pe
ventre	'e=rje	= 't ^h igi	=ã'ũn	= 'tæj	=rɛ'wɛk
sang	e='ruki	=ka' ⁿ bɾo	=d̄z=a'u	= 'hɛβp	=ru'i
brûler	= 'raĩ	=rɔ='k ^h ɹɔ̃	=po'k ^w a	mũ=...='haβp	=ra'pi
graisse	e='kira	= 't ^h wəmi	=d̄z=ap	= 'tuβp	= 'kap
pied	'e=i	= 'h ^w aji	= 'βi	=po'ta	= 'pi
main	'e=o	=nĩ'k ^h ɹa	= 'βo	= 'nĩβm	= 'pɔ
lourd	e='usi	=wi't ^h ĩ	=po'ti	=βp'təj	=pɔ'ij
foie	'e=ja	= 'nba	=pi'a	=tɛiβpkĩ'nāj	=pi'ʔa
nouveau	e='jasu	= 'ndiwi	=pa'gop	= 'tiβp	=pia'u
racine	e='rao	=ja'ɾe	kup=kujɔ'pe	mĩβm=nĩβm=tɛa'tiə	=ra'pɔ
peau	'e=i	= 'k ^h Λ	= 'pe	= 'tɛaj	= 'pit
queue	e='rokoĩ	= 'nbi	=d̄z=o'k ^w aj	=nā:='kiβp	= 'raj
blanc	'e=ʃĩ	=ja'k ^h a	=d̄zi'ra	=βp'douɥ	= 'sĩŋ
aile	e='heo	=ja'ɾa	=pe'o	=nĩ'māuɥ	=pe'pɔ, =ji'wa

Ci-dessous sont deux arbres construits en utilisant les mêmes listes. Certaines données (noms de langues et distances lexicostatistiques) manquent. Remplissez les trous. Pour chacun des arbres, précisez s'il est manuel ou automatisé, ainsi que s'il a été généré en utilisant l'Algorithme A ou B.



A	B	Γ	Δ	E
?	?	?	?	?

⚠ Les annotations manuelles et les indices de stabilité ont été intentionnellement omis pour cette tâche.

(M) (10 points) Les processus automatisés basés sur les classes Dolgopolsky peuvent produire des résultats erronés. Dans cet exemple le processus automatisé détecte plus de similarités entre le siriono et une certaine langue macro-jê (khîsêtjê), qu'entre le siriono et les autres langues tupi. Proposez un processus automatisé modifié qui produirait le classement correct s'il devait être appliqué aux listes de mots des familles macro-jê et tupi en haut, et décrivez ce processus *brièvement*.

⚠ Cette tâche ne sera notée que dans le cas d'une égalité de points entre les meilleures équipes.

Les auteurs remercient Alejandra Vidal, Maria Konoshenko, Ilya Gruntov et Jamthô Suyá d'avoir répondu à leurs questions sur des langues spécifiques. —*Andreï Nikouline, Miléna Vénéva*

Rédacteurs : Ivan Derjanski (éditeur technique), Hugh Dobbs, Stanislav Gurévitch, Boris Iomdin, Liam McKnight, Andreï Nikouline (éditeur en chef), Aleksejs Peguševs, Jan Petr, Alexandre Piperski, Maria Rubinstein, Miléna Vénéva, Elysia Warner.

Texte français : Elysia Warner.

Bon courage !