

第二十一届国际语言学奥林匹克竞赛

巴西·巴西利亚·2024年7月23—31日

团队赛题目

词汇统计法是一组根据词汇的相似程度来估算语言之间亲缘关系远近的方法。这类方法通常需要一个由专家手动标注过的冗长的单词表，标注的内容是单词表中的哪些单词被认为是同源的。不过，语言学家有时也将这类方法用在通过特定的流程自动生成标注的单词表上。其中一种流程基于辅音类这一概念，由俄裔以色列籍语言学家阿哈龙·多尔戈波尔斯基于1964年引入。

P. p b ɸ β f v	K. k g x γ q ɕ χ w	Y. j ç (在词根的开头)	M. m ɱ
T. t d ɖ θ ð ʈ ɟ	R. r r̥ ɽ l ʎ ʝ ʎ ʎ	W. w ɱ (在词根的开头)	N. n ɲ ŋ ŋ
S. s z ʃ ʒ ʂ ʐ ʑ ʒ			Q. ʈ ɟ
H. h ʕ ɦ ʔ ʕ ɦ ʔ, 所有的元音, 以及不在词根开头的 j ç w ɱ			

多尔戈波尔斯基的辅音类

本题将展示若干个已标注的单词表，每个单词表对应着一个语系。标注通过下标的数字表示。基于这几个单词表，语言学家通过名为“StarlingNJ”的算法的两个简化版本画出了语言谱系树状图，并计算了每个单词的稳定指数。在每个树状图表格中，上方的树状图和稳定指数基于的是手动标注的单词表，而下方的树状图和稳定指数基于的是自动标注的单词表。每种标注方式又各自对应两个树状图，分别通过算法A和算法B得到。需要注意的是，在某些情况下，一个单词表有不止一个可能的树状图。如果遇到这种情况，本题会随机展示这些树状图中的一个。在树状图的每个节点上都标有一个词汇统计距离。该距离越大，语言之间的亲缘关系越近。因此，“词汇统计距离”更准确的说法应该是“反词汇统计距离”。为简洁起见，本题依然采用“词汇统计距离”这个术语。

稳定指数和词汇统计距离均四舍五入至小数点后两位。若小数点后第三位小于5，则向下取整，否则向上取整。例如：2.836四舍五入为2.84；0.705四舍五入为0.71；0.703四舍五入为0.70。只有向人类读者展示的数值会被四舍五入。换言之，运行算法的计算机“看到”的是四舍五入之前的数值。

请注意，某些单词已知或疑似是外来词。例如：卡迪维奥语单词 **joki** ‘盐’来自瓜拉尼语 **juki**；伊派语（梅萨格兰德方言）单词 **ʔa:ni** ‘年’来自西班牙语 **ano**。

在某些情况下，单词表对于单个含义会列出多个同义词，用逗号分隔。例如贝霍斯语中的‘脚’。

在以下语料中，所有前缀均用“=”符号分隔，所有后缀均用“-”符号分隔。有些单词必须加前缀。这些单词以“=”符号开头。

所有语料均采用国际音标转写。' = 主重音、_ = 次重音（比主重音弱）、ː = 长音、˘ = 超短音、XY = X和Y同时发音、ˊ = 高平调、ˋ = 低平调、ˆ = 降调、ʔ = 前声门化发音（发音前喉部短暂地阻塞气流）、ʕ = 外挤气音（发音时喉部短暂的阻塞气流）、◌◌ = 清音、◌◌ = 鼻音化发音（通过鼻腔发音）、◌◌ = 嘎裂声（低沉、沙哑的发音）、◌◌表示鼻腔在辅音发音前有气流通过、◌ᵏ = 送气辅音（发音时有一股气流喷出）、◌ᵂ = 唇音化辅音（发音时圆唇）、◌ = 硬腭化发音（发音时舌头的一部分靠近硬腭）。a、æ、ɛ、ɪ、i、ɔ、u、ə、ʌ、ɒ、ə、y、ø、∅为元音。其他特殊字符均为辅音。

△ 与本题涉及的语言有关的知识不会为解题带来任何优势。

第壹部分. 瓜伊库鲁语系（阿根廷、巴西、巴拉圭）

	托巴语（东部方言）	皮拉加语	莫科维语（查科方言）	卡迪维奥语
云	l=ʔok ₁	'lo=ʔok ₁	naweyelek ₂	lol:adi ₃
火	nodek ₁	'd=oleʔ ₂	norek ₁	n=ol:edi ₂
鱼	njaq ₁	'nijaq ₁	naʎin ₂	nij:ogo-dʒegi ₃
头	=qajk ₁	'qajk ₁	=qaik ₁	=ak:ilo ₂
杀	=alawat ₁	=a'la:t ₁	=alawat ₁	=el:owadi ₁
月亮	ʔawoʒok ₁	ʔa'woʒok ₁	ʃirajɣo ₂	ep:enaj ₃
鼻子	=mik ₁	'mik ₁	=mik ₁	=m:iq:o ₁
盐	towe ₁	ol'ɣek ₂	ʔwe ₁	jok:i ₁
石头	qaʔ ₁	'qaʔ ₁	qaʔ ₁	wet:iga ₂
舌头	=atʃ-aʒat ₁	=a'tʃ-aʒat ₁	=oʔley-aʒan-aʒat ₂	=ok:el:i ₃

	算法A	算法B	
手动	<p>词汇统计距离</p>		云 0.50 火 0.50 鱼 0.50 头 0.75 杀 1.00 月亮 0.50 鼻子 1.00 盐 0.67 石头 0.75 舌头 0.50 稳定指数:
自动			云 0.50 火 0.50 鱼 0.75 头 0.75 杀 1.00 月亮 0.50 鼻子 1.00 盐 0.25 石头 0.75 舌头 0.50 稳定指数:

第贰部分. 努比亚语系（埃及、苏丹）

	敦戈拉维语	克努兹语	迪灵语	卡达鲁语	德布里语	比尔吉德语
杀	'bɛ:₁	be:₁	hur₂	wur-i₂	wur-i₂	fila:l-e₁
月亮	u'n-at-t₁	an-at-ti₁	nən-ti₁	nən-tu₁	nən-to₁	ma:l₂
水	'ɛss₁	essi₁	ɔti₁	ɔto₁	ɔtu₁	eji₁
给	'tir₁	tir₁	ti₁	ti₁	ti₁	te:-n₁
好	'sɛrɛ:₁	sere:₁	ken₂	kɛn₂	kɛŋ₂	azze-n₃
风	'turug₁	turug₁	irʃ-i₂	irʃ-o₂	irʃ-o₂	kurr-i₃
头发	'dɪl-ti₁	si:r₂	tel-ti₁	til-tu₁	til-tu₁	ur=dill-e₁
肚子	'tu:₁	tu:₁	te-te₂	to₁	to₁	tu:₁
睡觉	'nɛ:r₁	ne:r₁	ɟer₁	dwallɛli₂	ɟer-i₁	ne:r-i₁
太阳	'masɪl₁	masɪl₁	ɛɟ₂	aju₂	ɛŋgal-to₃	ʔi:zi₂

	算法 A	算法 B	
手动			稳定指数: 杀 0.50 月亮 0.83 水 1.00 给 1.00 好 0.50 风 0.50 头发 0.83 肚子 0.83 睡觉 0.83 太阳 0.50
自动			稳定指数: 杀 0.33 月亮 0.50 水 0.50 给 0.67 好 0.50 风 0.50 头发 0.83 肚子 1.00 睡觉 0.50 太阳 0.50

- (A) (2分) 辅音 **ɟ** 音似法语的 *r*，发音部位在舌头后端。该辅音属于哪一个多尔戈波尔斯基类？你是如何得出这一结论的？
- (B) (2分) 左上方的努比亚语系树状图是该算法和标注方式下两种可能的树状图之一。画出另一个可能的树状图。
- (C) (2分) 左下方的努比亚语系树状图是该算法和标注方式下两种可能的树状图之一。画出另一个可能的树状图。
- (D) (2分) 正如本题中展示的其他词汇统计距离一样，右上方树状图根节点的距离 0.49 只保留了两位小数。该距离在四舍五入前的具体数值是多少？

第参部分. 马塔瓜亚语系（阿根廷、玻利维亚、巴拉圭）

	维芝语（贝 尔梅霍河下 游方言）	维芝语 （里瓦达 维亚方 言）	贝霍斯 语	文纳耶克 语	伊约瓦雅语	曼惠语	尼瓦克勒 语（下游 方言）	尼瓦克勒 语（上游方 言）	马卡语
火	ʔitox ₁	ʔitox ₁	ʔitah ₁	ʔi:tax ₁	ʰwat ₂	ʔeite ₁	ʔitax ₁	ʔitax ₁	feʔt ₂
鱼	ʔwihat ₁	wahat ₁	wahat ₁	ʔwa:hat ₁	siʔjus ₋₁	ʃiʔjus ₋₁	saxetʃ ₋₁	saxetʃ ₋₁	sehets ₋₁
脚	=patʃu ₁	=qolob ₂	=patʃo ₁ , =kala ₂	=pa:kʔo ₁	=ʔsat ₃	=kaʔla ₂	=φo ₄	=φo ₄	=fʔi ₅
水	ʔinot ₁	ʔinot ₁	wah ₂	ʔina:t ₁	ʔiʔnat ₁	ʔaʔnat ₁	jinaʔt ₁	jinaʔt ₁	iweli ₃
给	=ʔwen _{-u} ₁	=wen _{-u} ₁	=ʔwen _{-o} ₁	=ʔwen _{-o} ₁	=ʔwehn-aʔm ₂	=ʔhaj ₃ , =ʔwen ₂	=xut ₄	=xut-ej ₄	tis-ix ₅
好	ʔis ₁	ʔis ₁	ʔis ₁	ʔis ₁	ʔes ₁	ʔeis ₁	ʔis ₁	ʔis ₁	t=ejkʔun-ej ₂
风	ʔinwok ^w ₁	ʔinwok ₁	ʔihwok ^w ₁	=ja:ʔ ₂ , =x ^w ox ^w ₃	ʔhlahwu ₄	ʔhlahwu ₄	ʔaβiʔm ₅	ʔaβiʔm ₅	tʔunikʔi ₆
树	haʔlo ₁	halob ₁	haʔla ₁	haʔla ₁	ʔaʔla ₁	ʔaʔla-k ₁	ʔaʔkxi-juk ₂	jiʔkla ₁	naxka-k ₃
头发	=ʔwule-j ₁	=wule-j ₁	=ʔwole-j ₁	=ʔwo:le-ç ₁ , hi:lenax ₂	=ʔwole ₁	=ʔwole-j ₁	=ʔateʔtj ₃	=jeʔs ₄	=ʔewkux-its ₅
杀	=lon ₁	=lon ₁	=lan ₁	=la:ŋ ₁	=ʔlaʔan ₁	=ʔlan ₁	=kla ₁	=kla ₁	=lan ₁

	算法 A	算法 B	
手动			<p>稳定指数:</p> <ul style="list-style-type: none"> 火鱼 0.78 鱼脚 1.00 水 0.33 给 0.78 好 0.44 风 0.89 树 0.33 头发 0.78 杀 0.67 杀 1.00
自动			<p>稳定指数:</p> <ul style="list-style-type: none"> 火鱼 0.78 鱼脚 0.44 水 0.33 给 0.56 好 0.67 风 0.89 树 0.22 头发 0.67 杀 0.67 杀 1.00

第肆部分. 蒙古语系（中国、蒙古、俄罗斯）

(E) (10分) 仔细研究以下单词表。求与自动标注和手动标注对应的稳定指数。

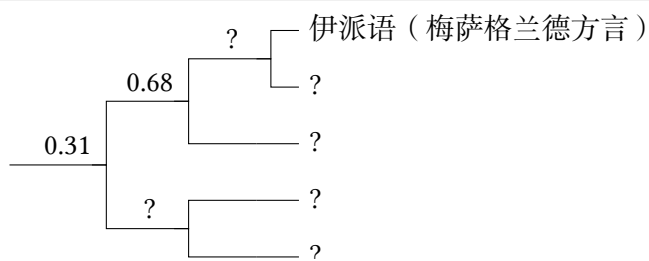
本题将单词‘所有’的两个稳定指数全部给出，以作提示。它们是：0.36、0.40（按随机顺序排列）。

	达斡尔语 (海拉尔方言)	哈密尼干语 (满洲方言)	布里亚特语 (豁里方言)	新巴尔虎语	额鲁特语	和硕特语	卡尔梅克语	喀尔喀语	鄂尔多斯语	东部裕固语	保安语
所有	hɔ:₁	bolt₂	buxi:₃	bygd₄	tsug₅	lug₅	tsuk₅, xamak₋₁	pux₃, pugt₄, xamăg₋₁	pyyite₄, xamukʰ₋₁	tʰuq₅	hanə₂
树皮	hails₁	qalihon₁	χoltɔhɔn₂	xalʰu:₁	xolts₂	xalis₁	dursn₃	xɔɣtʰɔs₂	turusu₃	χalsən₁	arasun₄
肚子	ke:li₁	gətəhən₂	gedehen₂	gedy:₂	ge:s₂	gets₂	gesn₂	gitis₂, xiwɣij₋₁	ketysy₂	ketesən₂	kele₁
鸟	dəgi₋₁	eiwan₁	ʃubu:n₁	ʃuwu:₁	ʃuvu:₁	ʃuwu:₁	ʃowun₁	ʃuwu₁	ʃuβu:₁	ʃu:n₁, peltʃɔr₂	bendzer₂
火	gali₁	gal₁	gal₁	gal₁	gal₁	gal₁	gal₁	gal₁	qal₁	qal₁	χal₁
路	terg-u:l₁	qargvi₂	χargi₂, zam₋₁	zam₋₁	dzam₋₁	dzam₋₁	xa:-lɔə₃	tsam₋₁	tʃam₋₁	mør₄	mor₄
盐	hata:₁	dawhon₂	dabhan₂	dawuhu:₂	daws₂	daws₂	dawsn₂	tawsă₂	taβusu₂	ta:psən₂	dabsuŋ₂
游泳	unpa-du₁	umba₋₁	tʰamar₋₂	umb₋₁	sele₋₃	umba₋₁	us-tɕi₋₄, ø:m₋₅	siɣi₋₃	usu-tʃi-la₋₄	umpa₋₁	mba₋₁
水	ɔsɔ₁	uxon₁	uhan₁	u:ha₁	usn₁	us₁	usn₁	ʊsɔ₁	usun₁	qʰusun₁	sə₁
风	kein₁	halkin₂	halxin₂	halxi₂	salʰxin₂	salkʲi₂	salʲkn₂	saɣxi₂	kʰi:₁	kʰi:₁	ki₁

第五部分. 尤玛语系（墨西哥、美国）

(F) (8分) 仔细研究以下单词表。下面是一个基于相同的单词表画出的树状图。其中缺少一些数据（语言名称和词汇统计距离）。填写其中的空缺。具体说明该树状图是基于手动标注还是自动标注，以及是用算法A还是算法B生成的。

	莫哈韦语	科科帕语	亚瓦派语	蒂派语（哈穆尔方言）	伊派语（梅萨格兰德方言）
短	wena=wen-a ₁	'xʌ=ʔut ₂	'tʃkr=ot-i ₂	lə=ʔuŋ ₁	mə=put-k ₃
鸟	ʔitʃ=i=jer ₁	'ʃa ₂	'ʔ=tʃ=sa ₂	aʔ='ʃa ₂	ʔa:=ʃa:₂
骨头	ŋ=a=s=ak ₁	'ŋ=j=a:k ₁	'tʃ=j=a:k-a ₁	'ak ₁	aq ₁
干燥	i=ro:-v-k ₁	's=ʔar ₂	'ru-β-i ₁	's=ʔa:j ₃	sa:j ₃
肉	kʷi:kway ₁	ʔi='ma:tʃ ₂	'kʷe:=ʔo-β-a ₃	'kʷak ₄	kukʷa:j-p ₁
脖子	maʎaqe ₁	'm=puk ₂	'mlq ₁	i:=ʔuk ₂	i:=ʔuk ₂
看到	i=ju:-k ₁	'wi:₂	'ʔu:₁	'wi:w ₂	ə=wu:w ₂
尾巴	i:=ʔar ₁	'ʃ=juʎ ₂	'β=hé ₃	ʃə='juʎ ₂	xə=juʎ ₂
二	havik-k ₁	'x=wak ₁	'hʷâk-i ₁	xə='wak ₁	xə=wak ₁
年	hu:ðe ₁	'mat-'ka:m ₂	'ʔ=tʃ ^h ur-a ₃	mat-'wam ₂	ʔa:n ^h ₁



(G) (20分) 语言学家还另外生成了一些尤玛语系树状图。这些树状图根节点的词汇统计距离（即树状图上最左边的词汇统计距离）如下：

1. 0.20
2. 0.23
3. 0.24

画出这些树状图。具体说明每一个树状图是基于手动标注还是自动标注，以及是用算法A还是算法B生成的。

(H) (3分) 小题(G)里词汇统计距离中的两个是在四舍五入至小数点后两位得到的：其中，0.23从0.225四舍五入得到。另一个在四舍五入后得到的距离是哪一个？该距离在四舍五入前的具体数值是多少？

(I) (4分) 写出稳定指数的计算方法。

(J) (5分) 写出词汇统计距离的计算方法。

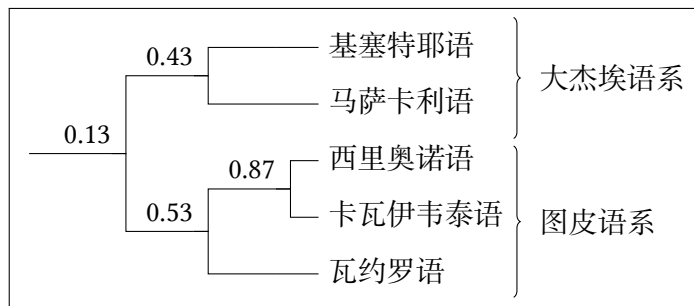
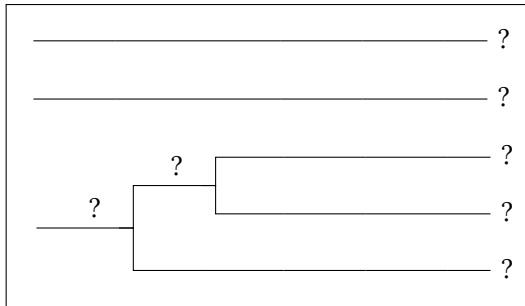
(K) (4分) 写出算法A与B之间的区别。

第陆部分.大杰埃语系、图皮语系（巴西、玻利维亚）

(L) (28分) 大杰埃语系和图皮语系是南美洲的两个主要语系。一些语言学家认为这两个语系存在远亲关系。仔细研究以下单词表。

	A	B	Γ	Δ	E
树皮	e='e-ke	h ^w i='k ^h Λ	kup='pε	mīβm='tεaj	= 'pε
肚子	'e=rje	= 't ^h igi	=ã'ũn	= 'tæj	=rε'wek
血	e='ruki	=ka' ⁿ bɾo	=d̄z=a'a	= 'hεβp	=ru'i
烧	= 'raī	=rɔ='k ^h Λs̄	=po'k ^w a	mũ=...='haβp	=ra'pi
肥	e='kira	= 't ^h wəmi	= 'd̄z=ap	= 'tuβp	= 'kap
脚	'e=i	= 'h ^w aji	= 'βi	=pv'ta	= 'pi
手	'e=o	=ɾi'k ^h Λa	= 'βo	= 'ɾiβm	= 'pɔ
重	e='usi	=wi't ^h i	=po'ti	=βp'təj	=pɔ'ij
肝	'e=ja	= 'nba	=pi'a	=tεiβpkā'nāj	=pi'ʔa
新	e='jasu	= 'ndiwi	=pa'gop	= 'tiβp	=pia'u
根	e='rao	=ja'ɾe	kup=kujɔ'pε	mīβm=ɾiβm=tεa'tiə	=ra'pɔ
皮肤	'e=i	= 'k ^h Λ	= 'pε	= 'tεaj	= 'pit
尾巴	e='rokoī	= 'nbi	=d̄z=o'k ^w aj	=nā:='kiβp	= 'raj
白	'e=fi	=ja'k ^h a	=d̄zi'ra	=βp'dou	= 'sɪŋ
翅膀	e='heo	=ja'ɾa	=pε'o	=ɾi'māu	=pε'pɔ, =ji'wa

下面是两个基于相同的单词表画出的树状图。其中缺少一些数据（语言名称和词汇统计距离）。填写其中的空缺。具体说明每一个树状图是基于手动标注还是自动标注，以及是用算法A还是算法B生成的。



A	B	Γ	Δ	E
?	?	?	?	?

⚠ 本小题省略了手动标注和稳定指数。

- (M) (10分) 基于多尔戈波爾斯基类的自动化流程可能产生错误的结果。在本题展示的例子中，该自动化流程的结果显示，西里奥诺语与大杰埃语系某个语言（基塞特耶语）的相似程度反而比西里奥诺语与图皮语系其他语言的相似程度高。提出一种优化的自动化流程，使得该流程能基于以上大杰埃和图皮语系单词表给出正确的分类结果，并对该流程进行简要的描述。

△ 本小题只有在高分团队平时时才计分。

作者感谢阿莱汉德拉·比达尔、玛丽娅·科诺申科、伊利亚·格伦托夫、亚姆托·苏亚回答有关特定语言的问题。
——安德雷·尼库林、米莱娜·韦内娃

编者：扬·彼得、戴谊凡（伊万·德尔然斯基）（技术编辑）、窦修（休·多布斯）、斯塔尼斯拉夫·古列维奇、玛丽娅·鲁宾斯坦、利亚姆·麦克奈特、安德雷·尼库林（主编）、阿莱克塞·佩古舍夫、亚历山大·皮佩尔斯基、王伊琳（埃莉西娅·沃纳）、米莱娜·韦内娃、鲍里斯·伊奥姆丁。

简体中文文本：刘羽扬。

祝你好运！