

Задача 3.

Съпоставителните изследвания в лингвистиката често се основават на *паралелни корпуси* — версии на едни и същи текстове на два (или повече) езика, които са *подравнени*, т. е. определено е и е отбелязано кои части отговарят една на друга. Подравняването обикновено се прави с помощта на програми, които обработват бързо големи количества текст, но нерядко допускат грешки, защото се ориентират само по формални критерии.

Ще разгледаме една такава програма, която приема два текста, разделени на редове, и разполага редовете един до друг, като дава на всяка двойка оценка, която е толкова по-висока, колкото по-уверена е програмата в съответствието. При това може някой ред от единия текст да обединява два или повече от другия или да не отговаря на нищо; не може обаче да се пренареждат редове. Програмата търси такава съотнасяне на редовете, при което средната оценка за целия текст да е колкото може по-висока.

Обикновено на програмата се дава абзац или по-дълъг текст, разделен на изречения, за да реши тя кои изречения отговарят едно на друго, но ние ще ѝ даваме по една дума на ред. Например да подравним българските изречения *В заведението влезе Марта с пазарска чанта в ръка* (П. Вежинов, „Везни“) и *Всичката тази плячка Авакум мигновено напъха във вътрешните джобове на сакото си* (А. Гуляшки, „История с кучета“) с украинските им преводи — съответно *До корчми зайшла Марта з господарською сумкою в руці* и *Всю цю здобич Авакум миттю порозпихав по внутрішніх кишенях піджака*:

<i>В</i>	<i>До</i>	0.1
<i>заведението</i>	<i>корчми</i>	0.0692308
<i>влезе</i>	<i>зайшла</i>	0.245455
<i>Марта</i>	<i>Марта</i>	1.8
<i>с</i>	<i>з</i>	0.3
<i>пазарска</i>	<i>господарською</i>	0.123529
<i>чанта</i>	<i>сумкою</i>	0.245455
<i>в</i>	<i>в</i>	1.8
<i>ръка</i>	<i>руці</i>	0.3

<i>Всичката</i>	<i>Всю</i>	-0.128571
<i>тази</i>	<i>цю</i>	0.06
<i>плячка</i>	<i>здобич</i>	0.3
<i>Авакум</i>	<i>Авакум</i>	1.8
<i>мигновено</i>	<i>миттю</i>	0.0818182
<i>напъха</i>	<i>порозпихав</i>	0.115385
<i>във</i>	<i>по</i>	0.18
<i>вътрешните</i>	<i>внутрішніх</i>	0.3
<i>джобове</i>	<i>кишенях</i>	0.3
<i>на ~~~ самото</i>	<i>піджака</i>	-0.0230769
<i>си</i>		-0.3

Както виждаме, програмата е разбрала, че двете български думи *на самото* отговарят на една украинска *піджака* (знакът ~~~ показва, че са обединени два реда от текста), а местоимението *си* няма съответствие.

а) В подравняването на изречението *И през това време Лоренцо наглеждаше* (А. Гуляшки, „Открадването на „Даная“) и превода му *І в цей час Лоренцо пильнував* са пропуснати оценките на съответствията:

<i>И</i>	<i>І</i>	
<i>през</i>	<i>в</i>	
<i>това</i>	<i>цей</i>	
<i>време</i>	<i>час</i>	
<i>Лоренцо</i>	<i>Лоренцо</i>	
<i>наглеждаше</i>	<i>пильнував</i>	

Във възходящ ред те са: $\{-0.3, 0.128571, 0.214286, 0.268421, 0.3, 1.8\}$.

Разположете ги на правилните места.

Решение:

Най-висока оценка (1.8) получава двойката, в която двете думи са еднакви. Оценките на останалите двойки зависят само от близостта на дължините на думите. От тях най-висока оценка (0.3) получава двойката, в която дължините са равни. От данните се вижда, че оценката е по-ниска, ако разликата в брой букви е по-голяма (двойката *това:цей* е по-добра както от *време:час*, така и от *през:в*), а при еднаква разлика тя тежи по-малко, ако думите са по-дълги (затова *наглеждаше:пильнував* е по-добра от *това:цей*). Остава да разпределим двете най-ниски оценки между *време:час* и *през:в*, а това можем да направим, като забележим, че *тази:цю* има оценка 0.06, а тази двойка е по-добра от *през:в*.

<i>И</i>	<i>И</i>	0.3
<i>през</i>	<i>в</i>	-0.3
<i>това</i>	<i>цей</i>	0.214286
<i>време</i>	<i>час</i>	0.128571
<i>Лоренцо</i>	<i>Лоренцо</i>	1.8
<i>наглеждаше</i>	<i>пильнував</i>	0.268421

б) Всъщност нищо не пречи вместо изречения, разделени на думи, да дадем на програмата просто списъци от думи. Разбира се, това няма да е съвсем коректно, ако знаем, че те отговарят една на друга. Но нека да видим какво става.

Ако двата „текста“ съдържат имената на месеците от годината съответно на български (1) и на украински (2), програмата предлага като най-добро следното подравняване (3):

(1)

януари
февруари
март
април
май
юни
юли
август
септември
октомври
ноември
декември

(2)

січень
лютий
березень
квітень
травень
червень
липень
серпень
вересень
жовтень
листопад
грудень

(3)

януари	січень	0.3
	лютий	-0.3
февруари	березень	0.3
март	квітень	0.1
април	травень	0.190909
май ~~~ юни	червень	0
юли	липень	0.0428571
август	серпень	0.253846
септември	вересень	0.264706
октомври	жовтень	0.26
ноември	листопад	0.26
декември	грудень	0.26

Както виждате, програмата е сметнала за вероятно, че украинското *лютий* няма българско съответствие, а *червень* отговаря на *май ~~~ юни*. Обединяването на два реда е намалило оценката за реда (в случая до 0), а полето без съответствие получава даже отрицателна оценка, но сумата за петте реда (0.290909) все пак е по-висока, отколкото ако не беше станало това изместване (тогава щеше да е 0.2748921):

февруари	лютий	0.136364
март	березень	0.0333333
април	квітень	0.190909
май	травень	-0.0428571
юни	червень	-0.0428571

Ако решим годината да започва от април, програмата предлага друго подравняване, като оставя този път без българско съответствие украинското *березень*:

(1)

<i>април</i>
<i>май</i>
<i>юни</i>
<i>юли</i>
<i>август</i>
<i>септември</i>
<i>октомври</i>
<i>ноември</i>
<i>декември</i>
<i>януари</i>
<i>февруари</i>
<i>март</i>

(2)

<i>квітень</i>
<i>травень</i>
<i>червень</i>
<i>липень</i>
<i>серпень</i>
<i>вересень</i>
<i>жовтень</i>
<i>листопад</i>
<i>грудень</i>
<i>січень</i>
<i>лютий</i>
<i>березень</i>

(3)

<i>април</i>	<i>квітень</i>	0.190909
<i>май ~~~ юни</i>	<i>травень</i>	0
<i>юли</i>	<i>червень</i>	-0.0428571
<i>август</i>	<i>липень</i>	0.3
<i>септември</i>	<i>серпень</i>	0.22
<i>октомври</i>	<i>вересень</i>	0.3
<i>ноември</i>	<i>жовтень</i>	0.3
<i>декември</i>	<i>листопад</i>	0.3
<i>януари</i>	<i>грудень</i>	0.253846
<i>февруари</i>	<i>січень</i>	0.207692
<i>март</i>	<i>лютий</i>	0.233333
	<i>березень</i>	-0.3

И има само един месец, с който, като започне годината, програмата дава правилното (от наше гледище) подравняване на месеците. Кой е той?

- | | | |
|-----------------------------------|---------------------------------|------------------------------------|
| <input type="checkbox"/> януари | <input type="checkbox"/> май | <input type="checkbox"/> септември |
| <input type="checkbox"/> февруари | <input type="checkbox"/> юни | <input type="checkbox"/> октомври |
| <input type="checkbox"/> март | <input type="checkbox"/> юли | <input type="checkbox"/> ноември |
| <input type="checkbox"/> април | <input type="checkbox"/> август | <input type="checkbox"/> декември |

Обяснете накратко отговора си.

Решение:

март

Разполагаме с две подравнявания, за които знаем, че програмата смята за по-добри от каноничното:

<i>януари</i>	<i>січень</i>	<i>април</i>	<i>квічень</i>
	<i>лютий</i>	<i>май ~~~ юни</i>	<i>травень</i>
<i>февруари</i>	<i>березень</i>	<i>юли</i>	<i>червень</i>
<i>март</i>	<i>квічень</i>	<i>август</i>	<i>липень</i>
<i>април</i>	<i>травень</i>	<i>септември</i>	<i>серпень</i>
<i>май ~~~ юни</i>	<i>червень</i>	<i>октомври</i>	<i>вересень</i>
<i>юли</i>	<i>липень</i>	<i>ноември</i>	<i>жовтень</i>
<i>август</i>	<i>серпень</i>	<i>декември</i>	<i>листопад</i>
<i>септември</i>	<i>вересень</i>	<i>януари</i>	<i>грудень</i>
<i>октомври</i>	<i>жовтень</i>	<i>февруари</i>	<i>січень</i>
<i>ноември</i>	<i>листопад</i>	<i>март</i>	<i>лютий</i>
<i>декември</i>	<i>грудень</i>		<i>березень</i>

Следователно верният отговор по условие не може да е *януари:січень* или *април:квічень*.

Ако предложим годината да започва с *февруари:лютий* или с който и да е месец от втората половина (VII – XII), програмата ще може да построи също такава подравняване като първото, само че с циклично изместване на месеците, което не променя общата оценка.

Ако поискаме годината да започва с *май:травень*, програмата ще може да отговори с второто подравняване с *април:квічень* преместени в края.

Ако годината започва с *юни:червень*, пак ще може да се получи по-добро според програмата подравняване, основано на второто: *април:квічень* и *май:травень* ще са в края,



https://dcl.bas.bg/cl_competition/

вместо *май* и *юни* (подравнени с *травень*) ще се обединят *юни* и *юли* (подравнени с *червень*), а на оценката това няма да се отрази, защото в *травень* и *червень* има еднакъв брой букви.

Така остава само *март:березень*. Не е лесно да се убедим, че при такова начало на годината няма възможно обединяване или пропускане на месеци, което да повиши оценката, но не е и нужно: достатъчно е (но това пък беше необходимо), че изключихме останалите 11 месеца.