

**МЕРЫ ЛЕКСИЧЕСКОГО СХОДСТВА
ЧАСТОТНЫХ СЛОВАРЕЙ**

**MEASURES OF LEXICAL SIMILARITY BETWEEN
FREQUENCY DICTIONARIES**

Аннотация. Предлагаются две меры лексического сходства частотных словарей (ЧС). Первая мера учитывает общее лексическое покрытие популяции двумя частотными словарями. Наибольший вклад в результат вносят частые слова. Вторая мера применима только к ЧС подкорпусов, она опирается на число общих лексических маркеров и в меньшей мере зависит от самых частых слов. Два подхода иллюстрируются данными ЧС русского, английского и китайского языков.

Ключевые слова. Лингвостатистика, подкорпусы.

Abstract. Two measures of lexical similarity between frequency dictionaries (FDs) are proposed. The first statistic measures common coverage of the population by two FDs. The main contribution to the result is made by most frequent words. The second measure is applicable to subcorpora of some general corpus. This statistic is based on the number of common keywords of FDs and to a lesser extent depends on most frequent words. The examples are drawn from FDs of Russian, English and Chinese.

Keywords. Statistical linguistics, subcorpora.

В лингвостатистике иногда возникает необходимость сравнить лексику двух частотных словарей и количественно определить их лексическое сходство.

В качестве простейшей меры такого сходства будем использовать формулу

$$C_{xy} = \sum \min \{p_{x_i}, p_{y_i}\} \quad (1)$$

Этот показатель приобретает значение 1 при сравнении частотных словарей (ЧС) одного и того же текста; он близок к 0 при сравнении ЧС, использующих разную графику.

Само собой разумеется, при работе с ЧС вместо вероятностей используются относительные частоты.

Можно заранее ожидать, что наибольший вклад в общую сумму внесут самые частые слова. Чтобы получить более дифференцированные результаты модифицируем формулу 1 следующим образом:

$$C_{xy} = \Sigma \min \{px_i, py_i\} / 0,5 (\Sigma px_i + \Sigma py_i), \quad (2)$$

где Σpx_i и Σpy_i подсчитываются для определенной зоны рангового словаря.

Покажем расчет на примере пяти самых частых слов из Частотного словаря современного русского языка, построенного на Национальном корпусе русского языка [Ляшевская, Шаров 2009: 419]. Сравним относительные частоты слов в Художественной литературе (ХЛ) и в Публицистике (Пуб.):

	ХЛ	Пуб.	min		ХЛ	Пуб.	min
и	.0365	.0355	.0355	на	.0167	.0167	.0167
в	.0262	.0364	.0262	я	.0183	.0104	.0104
не	.0223	.0169	.0169	Σ	.1200	.1159	.1057
C = .1057/.1180 = .896							

Подсчитаем коэффициент лексического сходства (C) для жанров этого же словаря, добавив к ХЛ и Пуб. еще два типа текстов: Устная речь (УР) и Прочее. Результаты приведены в табл. 1. В правой верхней части матрицы даются значения C для ста самых частых слов рангового словаря, в левой нижней части – для следующих четырехсот слов.

Таблица 1. Матрица лексического сходства жанров НКРЯ

	УР	ХЛ	Пуб.	Проч.
УР		.751	.667	.670
ХЛ	.661		.892	.852
Пуб	.561	.716		.933
Проч	.527	.700	.859	

Публицистика и прочая литература близки друг к другу. От них сильно отличается Художественная литература, совсем далека от них Устная речь.

Еще в 1960-е годы был создан Brown University Corpus of American English со следующими жанровыми подкорпусами (указана доля жанра в общем объеме – 1 миллион словоупотреблений):

A. Press: Reportage	8.8%	K. Fiction: General	5.8%
B. Press: Editorial	5.4%	L. Fiction: Mystery	
C. Press: Reviews	3.4%	and Detective	4.8%
D. Religion	3.4%	M. Fiction: Science	1.2%
E. Skills and Hobbies	7.2%	N. Fiction: Adventure	
F. Popular Lore	9.6%	and Western	5.8%
G. Belles Lettres, Biography, etc.	15.0%	P. Fiction: Romance and Love Story	5.8%
H. Miscellaneous	6.0%	R. Humor	1.8%
J. Learned and scientific Writing	16.0%		

[Kučerová, Francis 1967: 277-293, 100 самых частых слов]

Точно такой же объем и расклад жанров находим в Ланкастерском корпусе китайского языка (Lancaster Corpus of Mandarin Chinese – Oxford Text Archive), только жанр N сменил ярлык, здесь он называется Martial Art Fiction.

Таблица 2. Матрица лексического сходства жанров корпуса университета Брауна

	A	B	C	D	E	F	G	H	J	K	L	M	N	P	R
A		8 9	8 8	8 5	8 7	9 0	8 7	8 4	8 7	7 8	7 6	7 8	7 7	7 4	8 2
B			9 0	9 2	9 0	9 2	9 2	8 7	9 0	7 8	7 5	8 1	7 5	7 4	8 5
C				8 8	8 7	9 1	9 0	8 4	8 8	7 8	7 5	7 8	7 6	7 4	8 6
D					8 7	9 0	9 2	8 6	9 0	7 8	7 4	8 0	7 5	7 6	8 2
E						9 0	8 7	8 8	8 8	7 7	7 3	7 7	7 5	7 3	8 1
F							9 4	8 6	9 0	8 1	7 8	8 1	8 0	7 8	8 5
G								8 5	8 9	8 2	7 9	8 7	8 0	7 8	8 6
H									9 0	7 0	6 6	7 1	6 8	6 5	7 5
J										7 8	6 8	7 5	7 5	6 9	7 7
K											9 2	8 9	9 4	9 1	8 7
L												8 8	9 2	9 2	8 7
M													8 7	8 7	8 8
N														9 0	8 8
P															8 7
R															

Сто самых частых слов (со средним значением $C = .827$), ясно указывают на существование двух жанровых кластеров – ДЕЛОВОГО (А, В, С, D, E, F, G, J) со значениями C от .87 до .92 и ЛИТЕРАТУРНОГО (K, L, M, N, P) со значениями C от .87 до .94.

Результаты для китайского корпуса приведены в табл. 3. В правой верхней части матрицы даются значения C для ста самых частых слов рангового словаря, в левой нижней части – для следующих 1085 слов.

Среднее значение C в первой сотне китайских слов равно 74,1 (значительно ниже, чем в английском языке), во второй зоне оно равно 49,2. ЛИТЕРАТУРНЫЙ кластер помимо жанров K, L, M, N, P включает также R. ДЕЛОВОЙ кластер охватывает жанры В, С, D, E, F, J. Жанры А (репортаж) и G (биографии и эссе) располагаются между этими двумя кластерами.

Таблица 3. Матрица лексического сходства жанров Ланкастерского корпуса

	A	B	C	D	E	F	G	H	J	K	L	M	N	P	R
A		8 2	7 8	7 7	7 7	8 5	8 7	6 7	7 6	7 6	7 6	7 8	7 5	7 6	7 0
B	6 2		8 8	8 2	8 6	8 6	7 9	7 0	8 6	6 9	6 7	7 6	6 2	7 3	6 2
C	4 9	6 6		8 3	7 5	7 8	7 3	7 2	8 6	6 2	6 8	7 0	6 1	6 4	5 7
D	5 0	5 1	4 2		8 0	8 5	7 7	6 5	8 6	6 8	7 3	7 4	6 8	6 9	6 1
E	5 5	5 0	3 8	5 5		8 3	7 7	6 5	7 9	6 8	7 3	7 1	6 8	6 8	6 1
F	6 8	6 0	4 5	5 8	6 6		8 8	6 4	8 0	7 6	8 2	8 7	7 5	7 8	7 0
G	6 9	5 3	4 1	5 4	5 3	6 9		5 9	7 2	8 6	9 1	8 6	8 3	8 6	7 9

	A	B	C	D	E	F	G	H	J	K	L	M	N	P	R
H	4	4	5	3	3	3	3		7	5	5	5	4	5	4
	1	9	1	3	4	6	2		0	0	5	3	8	0	7
J	5	6	5	5	5	6	4	4		6	6	7	6	7	5
	3	1	3	8	4	0	7	8		4	8	1	2	2	7
K	5	4	2	4	5	5	6	2	3		8	8	8	9	8
	9	1	8	4	0	8	8	2	5		7	4	2	1	3
L	6	4	3	4	5	6	7	2	4	6		8	8	8	8
	2	6	3	8	4	3	0	7	2	4		6	5	7	1
M	5	4	3	4	4	5	5	2	4	4	5		7	8	7
	0	3	3	5	6	4	6	6	1	9	8		8	8	7
N	5	3	2	4	4	5	6	2	3	6	6	5		8	8
	2	8	5	4	9	1	5	0	2	8	6	1		1	2
P	5	4	3	4	5	6	6	2	3	7	7	5	6		8
	8	2	0	4	1	1	9	2	7	4	0	7	5		2
R	5	3	2	3	4	5	5	1	3	6	6	4	5	6	
	1	6	5	7	5	1	9	9	0	4	0	8	8	2	

По-другому выглядит конфигурация основных жанров в корпусе русской прозы 1850–1870-х гг. Здесь представлено восемь жанров – роман, повесть, рассказ, очерк, историческая проза, воспоминания, драма и прочее. Межжанровые значения *C* представлены в табл. 4 (справа зона 1–100 рангового словаря, слева – зона 601–3600).

В обеих зонах драма противостоит остальным жанрам, которые расположились цепочкой именно в том порядке, что дается в табл. 4. В зоне 601–3600 начинает проглядывать не столько жанровое, сколько собственно лексическое своеобразие.

Таблица 4. Матрица лексического сходства основных жанров прозы 1850–1870-х гг.

	Дра- ма	Ро- ман	Пов	Расск	Очер к	Ис- тор.	Восп	Проч
Драма		77	80	79	76	74	72	71
Роман	68		97	93	92	90	88	86
По- весть	70	88		95	92	90	89	85
Рассказ	68	81	84		92	92	88	86
Очерк	63	80	79	80		91	90	92
Истор.	58	71	71	72	68		87	88
Вос- пом.	61	77	76	72	76	64		90
Проч.	54	70	69	69	74	64	74	

Предложенная мера лексического сходства кажется очень естественной, ею можно воспользоваться при ЧС самого разного происхождения. Однако, у этой меры есть один существенный недостаток (как, впрочем, и у любой другой попытки представить словарь в виде одной цифры). Недостаток этот – обезличенность. За единым показателем скрываются сотни и тысячи реальных различий.

В какой-то степени преодолеть этот недостаток могут меры, опирающиеся именно на различия в ЧС. Рассмотрим подход, при котором сравниваемые ЧС относятся к какому-то уже существующему (или специально конструируемому на этом случай) корпусу с числом подкорпусов более двух. Для каждого из подкорпусов можно статистическими средствами выделить лексические маркеры. Используем для этого формулу

$$S = (x - m - 1) / \sqrt{m}, \quad (3)$$

где x – частота слова в подкорпусе, m – математическое ожидание этой частоты в рамках данной нулевой гипотезы.

Выбрав некоторый порог (скажем, $S > 3$), будем считать лексическими маркерами все слова, чьи S превысили данный порог.

Получив набор лексических маркеров для каждого из подкорпусов, мы для пары сравниваемых подкорпусов (X и Y) строим четырехклеточную таблицу вида

	X	не- X
Y	a	b
не- Y	c	d ,

где a – число общих маркеров X и Y , b – число маркеров Y , отсутствующих в X , c – число маркеров X , отсутствующих в Y , d – число маркеров, отсутствующих как в X , так и в Y .

Воспользуемся простейшим коэффициентом связи, введенным Дж. Юлом еще в 1900 г. [Upton, Cook: 435]:

$$Q = (ad - bc)/(ad + bc), \quad (4)$$

принимающим значения от -1 до $+1$.

Коэффициенты связи, полученные по формулам (3) и (4), не прошли еще проверки на широком материале ЧС. Исходя из общих соображений, можно ожидать интересных результатов при разнообразии подкорпусов в рамках общего ЧС. Складывается впечатление, что в зоне тысячи самых частых слов этот подход себя оправдывает.

Новые лингвостатистические инструменты открывают перед исследователем и новые (хотя и неясные) перспективы.

Литература

1. *Ляшевская О.Н., Шаров С.А.* (2009), Частотный словарь современного русского языка на материале Национального корпуса русского языка. М.
2. *Kučera F.* (1967), *Computational analysis of presentday American English*, Providence.
3. *Upton G., Cook I.* (2014), *A Dictionary of Statistics*, OUP.

References

1. *Lyashevskaya O.N., Sharov S.A. (2009), Chastotnyj slovar' sovremennogo russkogo yazyka na materiale Natsional'nogo korpusa russkogo yazyka. [Frequency Dictionary of Contemporary Russian on the Material of the National Corpus of Russian Language]. Moscow.*
 2. *Kučera F. (1967), Computational analysis of presentday American English, Providence.*
 3. *Upton G., Cook I. (2014), A Dictionary of Statistics, OUP.*
-

Шайкевич Анатолий Янович

Институт русского языка им. В.В. Виноградова РАН
(Россия)

Shaikevich Anatole

The Vinogradov Institute of Russian Language of the Russian
Academy of Sciences (Russia)

E-mail: lingstat@yandex.ru