

МАТЕМАТИЧЕСКИЕ МОДЕЛИ РОСТА СЛОВАРЯ  
И ИНФОРМАЦИОННЫХ ПОТОКОВ

В.В.Нешиной

При решении различных задач прогнозирования необходимо знать характер зависимости между изучаемыми величинами, например, между числом названий книг и брошюр или заявок на изобретения и временем; числом разных дескрипторов и числом заиндексированных документов; числом разных информационных запросов и числом абоненто-запросов (т.е. с учетом их повторяемости); числом разных слов частотного словаря и объемом выборки и т.д.

Целью настоящей статьи является отыскание уравнений, описывающих различного рода кривые роста из областей информатики и математической лингвистики и решение на их основе некоторых задач.

В качестве моделей кривых роста могут быть использованы различные формулы. Во-первых, это уравнения, найденные для описания кривых роста новых событий (под "новым" понимается любое из  $n$  разных событий, составляющих полную группу, при первом его появлении от начала испытаний). Во-вторых, это обобщенные плотности и функции распределения, задающие системы непрерывных распределений (при условии снятия ограничений, наложенных на параметры кривых распределения). Все эти модели содержат не более четырех параметров. Использование общих моделей значительно облегчает поиск выравнивающей кривой, хотя и не освобождает исследователя от глубокого анализа статистических данных и знания свойств указанных кривых. Как правило, частная модель содержит меньше параметров, чем их имеется в общей модели. Рассмотрим четыре основные системы кривых роста.

1. Система I кривых роста задается формулой (см. Нешиной В.В., 1986):

$$y = \int_0^x (1 - e^{-x\rho(t)}) dt, \quad (1)$$

где  $y$  - математическое ожидание числа разных событий (разных слов), наступающих при  $x$  испытаниях (появляющихся в выборке объемом  $x$  словоупотреблений);  $\rho(t)$  - непрерывная плотность распределения, аппроксимирующая вероятности  $\rho_k$  ( $k=1,2,\dots$ ,

$n$ ) разных событий, составляющих полную группу.

При заданной плотности  $\rho(t)$  по формуле (I) может быть рассчитана кривая роста новых событий. Все такие кривые удовлетворяют условиям

$$\left. \begin{aligned} \frac{dy}{dx} &= 1 \quad \text{при } x \rightarrow 0 \\ \frac{dy}{dx} &\rightarrow 0 \quad \text{при } x \rightarrow \infty \end{aligned} \right\} \quad (2)$$

В качестве плотности распределения  $\rho(t)$  в формулу (I) может входить любая плотность, заданная на интервале  $0 < t \leq n$ , где число разных событий  $n$  может быть как конечным, так и бесконечным. Например, это могут быть обобщенные плотности

$$\rho(t) = N t^{\gamma-1} (1 - \alpha u t^\beta)^{\frac{1}{u}-1}, \quad (3)$$

$$\rho(y) = \frac{N}{y} (\ln y)^{\gamma-1} (1 - \alpha u \ln^\beta y)^{\frac{1}{u}-1}. \quad (4)$$

Метод нахождения оценок параметров  $\alpha, \beta, \gamma, u$  обобщенных плотностей (3), (4) изложен в работе автора (Нешитой В.В., 1985).

Недостатком этой системы кривых роста является то, что интеграл (I), как правило, не выражается конечным числом элементарных функций. Кроме того, необходимо знать закон распределения вероятностей разных событий, составляющих полную группу.

## 2. Система П (а, б) кривых роста.

Система Па кривых роста в общем виде задается уравнением

$$\frac{dy}{dx} = 1 - \bar{F}(y) = 1 - \bar{F}(x), \quad (5)$$

где  $\bar{F}(y), \bar{F}(x)$  - функции распределения вероятностей новых событий.

Пусть  $\bar{F}(y), \bar{F}(x)$  задаются формулами (Нешитой В.В., 1986)

$$\bar{F}(y) = 1 - (1 - \alpha u y)^{\frac{1}{u}}, \quad (6)$$

$$\bar{F}(x) = 1 - [1 - \alpha(u-1)x]^{\frac{1}{u-1}}. \quad (7)$$

Тогда система Па кривых роста на основании (5) и (6), (7) будет описываться уравнением

$$y = \frac{1}{\alpha u} \left[ 1 - (1 - \alpha(u-1)x)^{\frac{u}{u-1}} \right], \quad (8)$$

которое удовлетворяет условиям (2).

Кривые роста, заданные общим уравнением (8), разделяются на типы. К I типу относятся кривые при  $0 < u < \infty$ ; ко II типу - при  $u \rightarrow 0$ ; к III типу - при  $-\infty < u < 0$ .

Оценки параметров  $\alpha, u$  кривых I типа рассчитываются по формулам

$$\alpha = \frac{1}{x^2} \sum_{m \neq 1} m^2 y_m - \frac{1}{x}, \quad u = \frac{1}{\alpha n},$$

где  $x = \sum_{m \neq 1} m y_m$ ,  $n = \sum_{m \neq 0} y_m$ ,  $y_m$  - число разных событий, наступающих при  $x$  испытаниях ровно  $m$  раз;  $n$  - число разных событий, составляющих полную группу.

Кривая роста относится к I типу, если  $y_{m=0} > 0$ .

В случае пуассоновского процесса ( $u \rightarrow 1$ ) уравнение (8) примет более простой вид

$$y = \frac{1}{\alpha} \left( 1 - \frac{1}{e^{\alpha x}} \right).$$

При  $u \rightarrow 0$  из (8) имеем кривую роста II типа

$$y = \frac{1}{\alpha} \ln(1 + \alpha x),$$

параметр  $\alpha$  которой находится по методу итераций

$$\alpha_{i+1} = \frac{1}{y} \ln(1 + \alpha_i x).$$

В случае кривых роста III типа, а также I и II типов, оценки параметров  $\alpha, u$  могут быть найдены упрощенным методом по трем известным из опыта величинам:  $x, y, y_{m=1}$ , т.е. по объему выборки  $x$ , количеству наступивших разных событий  $y$  и количеству событий с частотой  $m=1$ . При этом тип кривой устанавливается с помощью формулы

$$\frac{x}{y} = \frac{\frac{x}{y_{m=1}} - 1}{\ln \frac{x}{y_{m=1}}},$$

справедливой для кривых II типа. Если эмпирическое отношение  $(x/y)^* = x/y$ , то кривая роста относится ко II типу.

При  $(x/y)^* < x/y$  - к I типу. При  $(x/y)^* > x/y$  - к III типу.

Система Пб кривых роста получается из системы Па введением замены  $x = \ln X$ ,  $y = \ln Y$ . Тогда

$$\ln Y = \frac{1}{2u} \left[ 1 - (1 - 2(u-1) \ln X)^{\frac{1}{u-1}} \right] \quad (9)$$

или в дифференциальной форме

$$\frac{dy}{dx} = \frac{Y}{X} (1 - 2u \ln Y)^{\frac{1}{u}} = \frac{Y}{X} \left[ 1 - 2(u-1) \ln X \right]^{\frac{1}{u-1}} \quad (10)$$

Дифференциальное уравнение (10) при различных значениях параметра  $u$  дает ряд формул для описания кривых роста новых событий (например, новых слов в тексте), однако удобными для практического использования являются те из них, которые позволяют в явном виде выражать параметр  $\lambda$  через переменные  $X, Y$ . Это формулы при  $u = 1/2$ ,  $u = -1$  (см. табл. I).

Чтобы иметь более широкий набор кривых роста (с равным интервалом  $\Delta u = 0,75$ ), на основании формул с параметрами  $u = 1/2$ ,  $u = -1$  были получены еще две формулы, которые соответствуют значениям  $u \approx 1,25$ ,  $u \approx -0,25$ . Они тоже приведены в табл. I.

Проверка кривых роста новых слов в текстах показала, что параметр  $\lambda$  не является постоянной величиной, однако эмпирические точки в системе координат  $(\ln X; \lambda)$  ложатся на прямую

$$\lambda = \lambda_0 + k \ln X \quad (11)$$

Параметры  $\lambda_0, k$  являются параметрами текста и легко находятся по графику.

Эти же формулы оказались пригодными и для описания кривой роста новых слов в случайной выборке, но параметры выборки не совпадают с параметрами текста ( $\lambda_{0B} < \lambda_{0T}$ ,  $k_B > k_T$ ).

Примером такой кривой является кривая роста числа разных дескрипторов  $Y = f(X)$ , где  $X = Dh$  - число использованных дескрипторов при индексировании  $D$  документов;  $h$  - среднее число дескрипторов в одном документе.

Система 16 кривых роста новых слов в выборке

Таблица 1

Параметр $\mu$	Уравнение кривой роста	$\alpha = \alpha_0 + \kappa \ln X$	Оценки параметров выборки $\alpha_0, \kappa$
$\mu \approx 1,25$	$Y = X^{1/X} \alpha^{1/2}$	$\alpha = \frac{2}{\ln X} \ln \frac{\ln X}{\ln Y}$	$\kappa = \frac{2}{\ln^2 X} \left( 1 - \ln \frac{\ln X}{\ln Y} - \frac{Y_{m-1} \ln X}{Y \ln Y} \right)$ $\alpha_0 = \alpha - \kappa \ln X$
$\mu = 0,5$	$Y = X^{1/(1 + \frac{\alpha}{2} \ln X)}$	$\alpha = \frac{2}{\ln X} \left( \frac{\ln X}{\ln Y} - 1 \right)$	$\kappa = \frac{2}{\ln^2 Y} \left[ \left( \frac{\ln Y}{\ln X} \right)^2 - \frac{Y_{m-1}}{Y} \right]$ $\alpha_0 = \alpha - \kappa \ln X$
$\mu \approx -0,25$	$Y = X^{1/\sqrt{1 + \alpha \ln X}}$	$\alpha = \frac{1}{\ln X} \left[ \left( \frac{\ln X}{\ln Y} \right)^2 - 1 \right]$	$\alpha_0 = \frac{2}{\ln X} \left[ \frac{Y_{m-1}}{Y} \left( \frac{\ln X}{\ln Y} \right)^3 - 1 \right]$ $\kappa = \frac{1}{\ln^2 Y} - \frac{1}{\ln^2 X} - \frac{\alpha_0}{\ln X}$
$\mu = -1$	$Y = e^{\frac{1}{2} \sqrt{1 + 2\alpha \ln X} - 1}$	$\alpha = \frac{2}{\ln Y} \left( \frac{\ln X}{\ln Y} - 1 \right)$	$\kappa = \frac{2}{\ln^2 Y} \left[ 1 - \frac{Y_{m-1}}{Y} \left( 2 \frac{\ln X}{\ln Y} - 1 \right) \right]$ $\alpha_0 = \alpha - \kappa \ln X$

Оценки параметров выборки могут быть найдены по трем величинам:  $X, Y, Y_{m=1}$ . При этом используется формула В.М.Калинина (1964, с. 247)

$$\frac{d^m y}{dx^m} = (-1)^{m+1} \frac{m!}{x^m} y_m,$$

которая при  $m = 1$  дает

$$\frac{dy}{dx} = \frac{y_{m=1}}{x}. \quad (12)$$

На основании (12) и соответствующих уравнений кривой роста с учетом равенства (11) были найдены оценки параметров выборки, которые приведены в табл. I.

Из приведенных в табл. I формул лишь две последние ( $\mu \approx -0,25$ ,  $\mu = -1$ ) удовлетворяют условиям (2), при этом от параметра  $K$  зависит наибольший объем словаря. Кривая роста с параметром  $\mu \approx -0,25$  при  $X \rightarrow \infty$  дает:  $Y_{max} = e^{1/K}$ . Для кривой с параметром  $\mu = -1$  величина  $Y_{max} = e^{1/2K}$ . При  $K \rightarrow 0$  в обоих случаях  $Y_{max} \rightarrow \infty$ . Следовательно, параметр  $K$  может служить показателем лексического разнообразия (богатства) текста. При равных значениях  $K$  большим лексическим разнообразием обладает кривая с меньшим значением параметра  $\alpha_0$ .

Первые две формулы ( $\mu \approx 1,25$ ,  $\mu = 0,5$ ) описывают кривые, которые вначале возрастают, затем убывают (при весьма больших значениях  $X$ ). Следовательно, их можно использовать в качестве моделей кривых роста новых слов при ограниченных значениях  $X$  ( $X < 10^6 + 10^7$ ).

Формулы системы Пб кривых роста содержат всего два параметра, оценки которых легко находятся графическим методом либо рассчитываются по трем величинам  $X, Y, Y_{m=1}$ , при этом не требуется знание закона распределения разных слов по частоте их употребления в текстах.

3. Система III кривых роста задается формулами

$$y = NF(t), \quad (13)$$

$$y = N[1 - F(t)], \quad (14)$$

где  $N$  – некоторый параметр, причем,  $N = y_{\max}$ . Свойства кривых роста этой системы полностью определяются свойствами функции распределения  $F(t)$ .

Формула (13) может описывать, например, количество разных статей по определенной теме, опубликованных в первых  $t$  журналах, при условии, что последние упорядочены по убыванию количества таких статей, а также количество заболеваний при эпидемиях за время  $t$  от начала эпидемии.

#### 4. Система IV (а, б, в) кривых роста

Система IV кривых роста строится на основе обобщенных распределений. Вводя другие обозначения переменных и освобождаясь от ограничений, накладываемых на кривые распределения, можем записать следующие уравнения для описания различного рода кривых роста

$$\text{IVa: } y = Nx^{\gamma-1}(1-\alpha x^{\beta})^{\frac{1}{\alpha}-1} \quad (15)$$

$$\text{IVб: } y = Ne^{\gamma x}(1-\alpha e^{\beta x})^{\frac{1}{\alpha}-1}; \quad (16)$$

$$\text{IVв: } \ln y = N(\ln x)^{\gamma-1}(1-\alpha \ln^{\beta} x)^{\frac{1}{\alpha}-1}. \quad (17)$$

Система IV кривых роста является наиболее широкой системой, включающей кривые самой разнообразной формы. Она включает также некоторые кривые, принадлежащие другим системам.

#### 4.1. Система IVа кривых роста новых слов

Проверка показала, что общая формула (15) может достаточно точно описывать кривые роста новых слов при  $\gamma=2$ ,  $N=1$ :

$$y = x(1-\alpha x^{\beta})^{\frac{1}{\alpha}-1}. \quad (18)$$

Наиболее подходящей оказалась формула при  $\alpha = -1$  ( $0 < \beta < 1/2$ )

$$y = \frac{x}{(1+\alpha x^{\beta})^2}, \quad (19)$$

которая применима при  $(10^3 \div 10^4) < x < (10^7 \div 10^8)$  словоупотреблений. Оценки параметров текста находятся из уравнения прямой

$$\ln\left(\sqrt{\frac{x}{y}} - 1\right) = \ln \alpha + \beta \ln x. \quad (20)$$

Оценки параметров выборки можно найти по трем величинам:  $x, y, y_{m=1}$

$$\beta = \frac{1 - \frac{y_{m=1}}{y}}{2(1 - \sqrt{\frac{y}{x}})} = \frac{\sqrt{\frac{x}{y}}(1 - \frac{y_{m=1}}{y})}{2(\sqrt{\frac{x}{y}} - 1)}, \quad (21)$$

$$\alpha = \frac{1}{x^\beta} (\sqrt{\frac{x}{y}} - 1). \quad (22)$$

#### 4.2. Система ГУВ кривых роста новых слов

Уравнение (17) хорошо описывает кривую роста новых слов при  $\beta = 2, N = 1, (-1 \leq \alpha < 0)$ , т.е. система ГУВ кривых роста задается общей формулой

$$\ln Y = \ln X (1 - \alpha \ln^{\beta} X)^{\frac{1}{\alpha} - 1} \quad (23)$$

Пусть  $\alpha \rightarrow 0$ . Тогда из (23) получим уравнение

$$\ln Y = \frac{\ln X}{e^{\alpha \ln^{\beta} X}} \quad (X, Y \geq 1), \quad (24)$$

которое может быть приведено к прямой

$$\ln \ln \frac{\ln X}{\ln Y} = \ln \alpha + \beta \ln \ln X. \quad (25)$$

Параметры  $\alpha, \beta$  связного текста находятся путем построения по опытным значениям  $X_i, Y_i$  графика зависимости (25).

Параметры выборки можно рассчитать по трем величинам:  $X, Y, Y_{m=1}$

$$\beta = \frac{1}{\ln \frac{\ln X}{\ln Y}} \left( 1 - \frac{Y_{m=1}}{Y} \frac{\ln X}{\ln Y} \right), \quad (26)$$

$$\alpha = \frac{1}{\ln^{\beta} X} \ln \frac{\ln X}{\ln Y}. \quad (27)$$

Пусть далее  $\alpha = -1$ . Тогда из (23) получим

$$\ln Y = \frac{\ln X}{(1 + \alpha \ln^{\beta} X)^2}, \quad (28)$$

откуда

$$\ln \left( \sqrt{\frac{\ln X}{\ln Y}} - 1 \right) = \ln \alpha + \beta \ln \ln X. \quad (29)$$



Оценки параметров выборки равны

$$\beta = \frac{1 - \frac{Y_{m=1}}{Y} \frac{\ln X}{\ln Y}}{2 \left( 1 - \sqrt{\frac{\ln Y}{\ln X}} \right)}, \quad (30)$$

$$\alpha = \frac{1}{\ln^{\beta} X} \left( \sqrt{\frac{\ln X}{\ln Y}} - 1 \right). \quad (31)$$

Достоинством формул (24), (28) является то, что они хорошо работают практически от начала координат до весьма больших значений  $X$  ( $X = 10^7 + 10^8$  словоупотреблений). Однако при дальнейшем увеличении  $X$  поведение этих кривых не соответствует поведению кривых роста новых событий.

Опыт показывает, что параметры выборки в приведенных выше формулах зависят как от типа текстов, на основе которых построена данная выборка, так и от ее объема. В связи с этим полученные формулы остаются справедливыми только в пределах заданной выборки, т.е. их нельзя использовать для экстраполяции кривой роста новых слов на выборку большего объема.

В то же время параметры текста не зависят от объема текста (при условии его лексической однородности), что позволяет прогнозировать объем словаря с ростом объема текста.

Проверим работу формул IV и V (а, в) кривых роста новых слов в выборке. Восстановим кривую  $y=f(x)$  на основе опытных данных по точным формулам, а также рассчитаем ее по приближенным формулам, при этом параметры выборки  $\alpha_0, \kappa$ , а также  $\alpha, \beta$  определим по трем величинам:  $X, Y, Y_{m=1}$ .

Для восстановления кривой роста новых слов в выборке воспользуемся формулой В.М.Калинина (1964, с.247)

$$y = Y - \sum_{m \geq 1} \left( 1 - \frac{x}{X} \right)^m Y_m, \quad (32)$$

где  $Y$  - число разных слов в выборке объемом  $X$  словоупотреблений;  $Y_m$  - число слов с частотой  $m$  в выборке  $X$ ;  $y$  - ожидаемое среднее число разных слов в подвыборке произвольного объема  $x$  ( $x < X$ ).

Формулой (32) удобно пользоваться при  $x/X \geq 0,1$ .

При  $x/X < 0,1$  целесообразно воспользоваться формулой (I), в которую вместо плотности  $\rho(t)$  следует подставить относительные частоты слов  $\rho_z^* = m_z^*/X$ , где  $z$  - ранг слова в

частотном словаре. Тогда формула (I) примет вид

$$y = \int_0^n (1 - e^{-x\rho^z}) dz. \quad (I')$$

Интегрирование осуществляется численным методом по формуле прямоугольников.

Таблица 2  
Параметры выборок, рассчитанные по данным двух частотных словарей (ЧС)

Параметры	Пб (табл. I)			ГVa	ГVв	
	$\mu \approx 1,25$	$\mu = 0,5$	$\mu = -0,25$	$\mu = -1$	$\mu \rightarrow 0$	$\mu = -1$
ЧС немецкого языка						
$\alpha_0$	0,02029	0,01622	0,00996			
$\kappa$	0,0007485	0,001287	0,002017			
$\alpha$				0,04223	0,005719	0,002363
$\beta$				0,3005	1,3741	1,46635
ЧС русского языка						
$\alpha_0$	-0,001633	-0,01485	-0,03344			
$\kappa$	0,002935	0,004306	0,00615			
$\alpha$				0,01451	0,001262	0,000466
$\beta$				0,4084	2,0418	2,18324

В таблице 3 (столбец 2) приведены результаты расчетов по точным формулам (I') и (32) на основе опытных данных частотного словаря немецкого языка и по приближенным формулам систем кривых роста Пб, ГVa, ГVв (столбцы 3-8). Объем выборки здесь равен  $X = 10910777$ , объем словаря  $Y = 258173$  и количество одноразовых слов  $Y_{m=1} = 126862$  (Meier H., 1964). Параметры выборки для каждой аппроксимирующей кривой приведены в табл. 2.

Аналогичные расчеты выполнены по данным "Частотного словаря русского языка" (под ред. Л.Н. Засориной, 1977). Здесь  $X = 1056382$ ,  $Y = 39268$ ,  $Y_{m=1} = 13379$ .

Из табл. 3 видно, что в первом случае все формулы достаточно точно описывают кривую роста новых слов в выборке, при этом наименьшую точность показала формула при  $\mu \approx 1,25$  (система Пб). Во втором случае менее точными оказались две формулы из той же системы кривых роста: при  $\mu = 0,5$ ,  $\mu \approx -0,25$ .

Таблица 3

Кривые роста новых слов, восстановленные по двум ЧС и рассчитанные по формулам систем кривых роста Пб, Пв, Пв

Объем выборки X	Объем словаря по ЧС Y	IIσ			IVα	IVβ	
		u≈1,25	u=0,5	u=-0,25	u=-1	u→0	u=-1
ЧС немецкого языка							
1091	625	600	619	646	603	626	636
2000	1019	979	1010	1054	1000	1019	1037
5000	2118	2004	2066	2156	2092	2082	2116
10911	3878	3609	3716	3866	3819	3787	3794
30000	-	7503	7699	7965	8009	7725	7828
100000	17800	17118	17462	17912	18207	17485	17663
300000	-	34735	35215	35815	36459	35220	35468
1091078	77870	75440	75952	76545	77420	75913	76178
3273233	140810	139063	139377	139670	140230	139293	139474
10910777	258173	258167	258232	258173	258173	258083	258177
ЧС русского языка							
1000	625	650	713	-	645	646	659
3000	1510	1534	1674	-	1571	1526	1556
10000	3637	3618	3889	4274	3790	3604	3667
30000	7384	7296	7682	8193	7670	7275	7374
105638	14815	14703	15095	15569	15215	14683	14791
211276	20700	20597	20890	21227	21028	20582	20666
528191	30495	30455	30544	30645	30600	30450	30477
1056382	39263	39268	39269	39268	39268	39269	39269

\* IIσ

$$u \approx 1,25: Y = X^{1/X^{d/2}}$$

$$u = 0,5: Y = X^{1/(1 + \frac{d}{2} \ln X)}$$

$$u \approx -0,25: Y = X^{1/\sqrt{1 + d \ln X}}$$

$$(d = d_0 + k \ln X)$$

IV

$$a) u = -1: Y = \frac{x}{(1 + dx^\beta)^2}$$

$$b) \begin{cases} u \rightarrow 0: \ln Y = \frac{\ln X}{e^{d \ln^\beta X}} \\ u = -1: \ln Y = \frac{\ln X}{(1 + d \ln^\beta X)^2} \end{cases}$$

Полученные результаты свидетельствуют о том, что параметр  $\alpha$  изменяется от выборки к выборке. А это значит, что для более точного описания кривой роста новых слов в выборке необходимо использовать общие формулы (9), (18), (23), каждая из которых (с учетом равенства (II) в случае формулы (9)) содержит три параметра. Для приближенных же расчетов удовлетворительные результаты дают рассмотренные выше формулы, имеющие всего по два параметра.

Возможность приведения этих формул к прямой оказалась весьма полезной для решения различных задач. Рассмотрим некоторые из них.

#### 4.3. Оценка степени аналитичности языка

Показателем степени аналитичности языка принято считать коэффициент, равный отношению количества лексем к количеству словоформ частотного списка (Пиотровский Р.Г. и др., 1962). Степень аналитичности языка тем выше, чем ближе этот показатель к единице.

Существенным недостатком этого показателя является его зависимость от объема текста. Знание аналитической зависимости между объемами словаря и текста позволяет по-другому измерять степень аналитичности языка.

Пусть кривая роста новых слов в тексте описывается формулой (19). Величина параметра  $\beta$  зависит от выбора единицы подсчета количества разных слов, в качестве которой может быть принята словоформа (СЛ) или лексема (Л). При этом для одного и того же текста  $\beta_L > \beta_{СЛ}$ , в то время как  $\alpha_L \approx \alpha_{СЛ}$ . Последнее равенство позволяет ввести показатель степени аналитичности языка, который не зависит от объема текста:

$$\Delta\beta_a = \frac{v_L - v_{СЛ}}{\ln x}, \quad (33)$$

где  $v_L, v_{СЛ}$  рассчитываются по формуле

$$v = \ln \left( \sqrt{\frac{x}{y}} - 1 \right). \quad (34)$$

Чем ближе  $\Delta\beta_a$  к нулю, тем выше степень аналитичности языка.

В табл.4 приведены значения  $\Delta\beta_a$  для четырех языков, рассчитанные по формулам (33), (34) для текстов по электронике.

Полученные результаты позволяют осуществить переход от кривой роста новых словоформ в тексте к кривой роста новых лексем. Для этого достаточно воспользоваться формулами

$$\alpha_L = \alpha_{СЛ}; \quad \beta_L = \beta_{СЛ} + \Delta\beta_a.$$

Таблица 4  
Степень аналитичности языков (тексты по электронике)

Язык, источник	Объем текста $x$	Объем словаря		$\Delta\beta_a$
		$Y_{сл}$	$Y_c$	
Русский (Калинина, Е.А., 1968)	200894	21468	6826	0,0627
Румынский (Ешан Л.И., 1966)	200000	14292	5708	0,0479
Французский (Кочеткова В.К., Скредина Л.М., 1968)	100000	8108	4527	0,0336
Английский (Алексеев П.М., 1968)	200000	10582	7160	0,0202

#### 4.4. Показатель степени связности слов в лексически однородном тексте

Так как в связанном тексте лексико-грамматические связи накладывают определенные ограничения на сочетаемость слов, то естественно предположить (и это подтверждается опытными данными), что число разных слов в отрезке сплошного текста в среднем будет меньше, чем в случайно составленной выборке равного объема, взятой из достаточно большой совокупности лексически однородных текстов.

Эту разницу в объемах словарей можно использовать для оценки степени связности слов в тексте, причем, она оказывается независимой от объема текста.

Пусть для некоторого целого произведения из опыта известны объемы - всего текста  $X$ , словаря  $Y$ , одноразовых слов  $Y_{m=1}$ , а также несколько промежуточных точек  $(x_i; y_i)$  на кривой роста новых слов. По этим данным найдем параметры текста  $\alpha_1, \beta_1$  путем построения графика зависимости (20). Для этого же произведения можно рассчитать параметры выборки  $\alpha_2, \beta_2$  по формулам (21), (22). Эти параметры будут относиться к такой кривой роста, которая получилась бы при случайном отборе словоупотреблений из данного произведения и подсчете количества разных слов.

Многочисленные расчеты показали, что параметры  $\alpha_1$  и  $\alpha_2$  находятся в отношении

$$\frac{\alpha_1}{\alpha_2} \approx 2, \quad (35)$$

которое может быть использовано в качестве показателя степени связности слов в тексте.

Параметры  $\beta_g$  и  $\beta_r$  связаны соотношением

$$\beta_r = \beta_g - \frac{\lg 2}{\lg X} \quad (36)$$

формулы (35), (36) позволяют по параметрам выборки (которые легко вычислить по трем величинам  $X, Y, Y_{m-1}$ ) найти параметры текста.

#### 4.5. Оценка лексической близости двух связанных текстов.

##### Автоматическая классификация текстов

Пусть кривая роста новых слов в связанном тексте описывается формулой (19). Объединим два равных по объему текста с одинаковыми параметрами  $\alpha, \beta$  и исследуем поведение обоих параметров этого объединенного текста. Рассмотрим два крайних случая.

Случай 1. Параметры  $\alpha, \beta$  объединенного текста равны соответствующим параметрам объединяемых текстов, а его словарь содержит такое же количество разных слов, сколько их было бы в каждом из объединяемых текстов при удвоенной его длине. В этом случае степень лексической близости двух текстов (будем измерять ее некоторым показателем  $\zeta$ ) равна единице, а точка с координатами  $(\ln x_{12}; v_{12, \zeta=1})$  лежит на прямой 1 (см. рис. 1), поскольку

$$v_{12, \zeta=1} = \ln \left( \sqrt{\frac{2x_1}{y_{12, \zeta=1}}} - 1 \right) = \ln \alpha + \beta \ln 2x_1 = v_1 + \beta \ln 2.$$

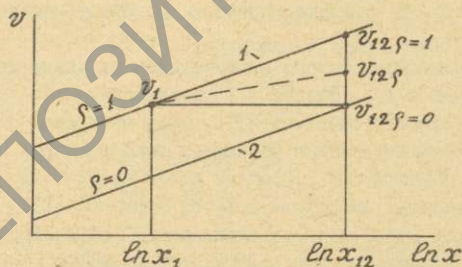


Рис. 1. Графическое определение показателя  $\zeta$

Случай 2. Словари обоих текстов не содержат общих слов. Это значит, что степень лексической близости таких текстов равна нулю ( $\zeta = 0$ ). В этом случае объем словаря объединенного текста равен удвоенному объему словаря одного из объединяемых текстов, а величина

$$v_{12\gamma=0} = \ln \left( \sqrt{\frac{2x_1}{2y_1}} - 1 \right) = v_1 = \ln d + \beta \ln x_1,$$

т.е. прямая 2 параллельна прямой 1 и расположена ниже ее на расстоянии  $\beta \ln 2$ .

Из рис. 1 видно, что величина  $v_{12\gamma}$ , характеризующая объем словаря реального объединенного текста, ограничена  $v_{12\gamma=0} < v_{12\gamma} < v_{12\gamma=1}$ . Следовательно, показатель  $\gamma$  можно измерять отношением

$$\gamma = \frac{v_{12\gamma} - v_{12\gamma=0}}{v_{12\gamma=1} - v_{12\gamma=0}} = \frac{v_{12\gamma} - v_{12\gamma=0}}{\beta \ln 2}. \quad (37)$$

Из условия параллельности прямых 1 и 2 следует, что  $\gamma$  не зависит от объема текстов  $x_i$ , но оба они должны быть равными между собой.

Отметим, что для двух сравниваемых текстов отношение  $y_{12\gamma=0}/y_{12\gamma=1}$  не может превосходить определенного значения, которое при  $0 < x_i < \infty$  заключено на интервале  $1 < y_{12\gamma=0}/y_{12\gamma=1} < 4^\beta$ .

Например, для "Частотного словаря русского языка" при  $\beta = 0,4084$  (табл. 2, столбец 5) это отношение не должно превышать величины  $4^\beta = 1,7615$ . Действительно, при  $x_1 = 528191$  по данным табл. 3 (столбец 2) имеем:  $30495 \cdot 2/39268 = 1,5532$ .

Параметр  $\beta$ , входящий в формулу (37), может быть надежно определен лишь по тексту достаточно большого объема. В случае коротких текстов показатель  $\gamma$  можно оценить проще, если объем объединенного текста принять равным  $x_{12} = x_1/2 + x_2/2$ . Тогда для вычисления показателя  $\gamma$  достаточно будет найти значения следующих величин:

$$y_{12\gamma=0} = y_1(x_1/2) + y_2(x_2/2); \quad y_{12\gamma} = y(x_{12});$$

$$y_{12\gamma=1} = \frac{1}{2} [y_1(x_1) + y_2(x_2)],$$

где  $y_1(x_1/2)$  - количество разных слов среди половины словоупотреблений первого текста (например, стоящих только на четных местах);  $y_2(x_2/2)$  - то же для второго текста;  $y(x_{12})$  - количество разных слов среди половины словоупотреблений первого и второго текстов;  $y_1(x_1)$ ,  $y_2(x_2)$  - количество разных слов соответственно в первом и втором текстах.

При небольших размерах текстов ( $x < 40000$ ) кривая роста новых слов достаточно точно описывается формулой

$$y = \frac{x}{1 + dx^\beta},$$

которая следует из (18) при  $u \rightarrow -\infty, \Delta u > 0$ . Ее можно привести к виду

$$v = \ln\left(\frac{x}{y} - 1\right) = \ln \Delta + \beta \ln x. \quad (38)$$

Для определения показателя  $\beta$  (см. формулу (37)) необходимо по формуле (38) вычислить значения  $v_{12\beta}$  при заданных  $x_{12}$  и  $y_{12\beta}$ . При этом получим оценку  $\beta$ , не зависящую от объема сравниваемых текстов (но оба текста должны быть равными между собой).

Все операции по вычислению  $\beta$  легко автоматизировать. Задавая пороговое значение этого показателя, можно классифицировать тексты (документы) по тематическим группам, включая в одну группу близкие по содержанию тексты.

#### 4.6. Система IVa кривых роста информационных потоков

Одной из таких кривых, принадлежащих данной системе, является кривая роста (во времени  $x$ ) числа названий книг и брошюр. Ее свойства отличаются от свойств кривых роста новых слов. При этом, как показывает анализ опытных данных, рост числа названий книг и брошюр, заявок на изобретения и т.д. может быть либо линейным, либо близким к экспоненте. Следовательно, общее уравнение для описания этих кривых роста может быть записано в виде (при  $\beta = 1, \Delta < 0, 0 < u < 1$ )

$$y = y_0 (1 - \Delta u x^\beta)^{\frac{1}{u} - 1}, \quad (39)$$

где  $y_0$  - начальное значение  $y$  при  $x = 0$ .

Для нахождения параметров  $\Delta, \beta, u$  преобразуем (39) к виду

$$\ln\left[\left(\frac{y}{y_0}\right)^{\frac{1}{1-u}} - 1\right] = \ln \Delta u + \beta \ln x. \quad (40)$$

Значение параметра  $u$  должно быть подобрано таким, чтобы график зависимости (40), построенный по опытным значениям  $y, x$ , представлял собой прямую. Оценки параметров  $\Delta, \beta$  могут быть определены по графику полученной прямой.

В случае, если ни при каких значениях параметра  $u$  ( $0 < u < 1$ ) прямая не получается, то это свидетельствует о том, что формула (39) в данном случае не работает и следует рассмотреть кривую другого типа, например,

$$y = y_0 e^{\Delta x^\beta}. \quad (41)$$



Оценки параметров  $\alpha, \beta$  находятся на основании уравнения прямой

$$\ln \ln \frac{y}{y_0} = \ln \alpha + \beta \ln x \quad (42)$$

при  $\alpha > 0$ , либо прямой

$$\ln \ln \frac{y_0}{y} = \ln \alpha + \beta \ln x \quad (43)$$

при  $\alpha < 0$ . Исследования показывают, что кривая роста (41) при  $\beta < 0$  имеет горизонтальную асимптоту  $y = y_0$ , которую также следует оценить по опытным данным. Для этого представим формулу (41) в виде

$$y_{nx} = \alpha y_x^b, \quad (44)$$

где  $y_{nx}$  может обозначать число наступивших разных событий в выборке объемом  $nx$  (за время  $nx$ ). Параметры  $\alpha, b$  находятся из уравнения прямой

$$\ln y_{nx} = \ln \alpha + b \ln y_x. \quad (45)$$

Тогда оценки параметров в формуле (41) будут равны

$$y_0 = \alpha^{1/(1-b)} \quad (b \neq 1), \quad (46)$$

$$\beta = \frac{\ln b}{\ln n}, \quad (47)$$

$$\alpha = \frac{1}{3} \sum_{i=1}^3 \frac{1}{x_i^\beta} \ln \frac{y_i}{y_0}. \quad (48)$$

Если  $b = 1$ , то кривая роста задается формулой  $y = \alpha x^\beta$ .

Формулы (41)–(45) хорошо описывают кривые роста числа названий книг и брошюр, кривые роста производительности труда и произведенного национального дохода (Нешитой В.В., 1984) и некоторые другие кривые. Важной характеристикой этих кривых является темп роста

$$q_i = \frac{y_i}{y_{i-1}} = e^{\alpha(x_i^\beta - x_{i-1}^\beta)}. \quad (49)$$

Из (49) следует, что при  $\beta = 1$   $q_i = e^\alpha = \text{const}$ , т.е. темп роста не зависит от времени  $x_i$  и равен темпу роста показательной функции  $y = e^{\alpha x}$ . При  $\beta > 1$  величина  $q_i$  с ростом  $x_i$  растет, а при  $\beta < 1$  — уменьшается, т.е. параметр  $\beta$  является показателем ускорения или замедления темпа роста кривой (41). При  $x_1 = 1$ ,

$x_0 = 0$  из (49) имеем  $q_1 = y_1/y_0 = e^{\alpha}$ , откуда  $\alpha = \ln(y_1/y_0) = \ln q_1$ , т.е. величина  $\alpha$  представляет собой натуральный логарифм темпа роста кривой на начальном отрезке времени, равном единице. Поскольку  $q_1 = e^{\alpha} = 1 + \frac{\alpha}{1!} + \frac{\alpha^2}{2!} + \dots$ , то при малых  $\alpha$  ( $\alpha < 0,1$ ) имеем  $q_1 \approx 1 + \alpha$ , откуда  $\alpha \approx q_1 - 1$ .

Проверим работу формул (41), (44) на практическом примере. Воспользовавшись опытными данными (Михайлов А.И. и др., 1976, с.212, 213), приведенными в табл.5, столбцы 1,3, найдем параметры выравнивающей кривой роста числа названий книг и брошюр, изданных в 16 странах за 1958-1972 г.г. Пусть эмпирическая кривая роста описывается формулой (41), которую представим в виде

$$\frac{y}{y_0} = e^{\alpha x^{\beta}}, \quad (50)$$

где  $y_0 = 216696$  (количество книг и брошюр на конец 1958 г.). Приведем ее к уравнению прямой

$$\left(\ln \frac{y}{y_0}\right)_{nx} = \beta \left(\ln \frac{y}{y_0}\right) x \quad (51)$$

Таблица 5

Число названий книг и брошюр ( $y$ ), изданных в 1958-1972 г.г. в 16 странах (Михайлов А.И. и др., 1976)

Год	$x$	$y$	$\frac{y}{y_0}$	$\ln \frac{y}{y_0}$	$\left(\frac{y}{y_0}\right)^{расч.}$
1958	0	216696	1	0	1
1959	1	224201	1,0346	0,0341	1,0378
1960	2	237837	1,0976	0,0931	1,0824
1961	3	245382	1,1324	0,1243	1,1314
1962	4	257619	1,1889	0,1730	1,1843
1963	5	269942	1,2457	0,2197	1,2410
1964	6	282892	1,3055	0,2666	1,3016
1965	7	288422	1,3310	0,2859	1,3662
1966	8	323388	1,4924	0,4004	1,4350
1967	9	333848	1,5406	0,4322	1,5081
1968	10	336975	1,5551	0,4415	1,5858
1969	11	347699	1,6045	0,4728	1,6682
1970	12	385570	1,7793	0,5762	1,7557
1971	13	390247	1,8009	0,5883	1,8486
1972	14	396661	1,8305	0,6046	1,9471

Если построить по эмпирическим значениям величин  $ln(y/y_0)$  график зависимости (51) при  $n = 2$ , т.е. при всех возможных парах значений  $(x, 2x)$ : 1,2; 2,4; 3,6;...; 7,14, - то убедимся, что точки действительно рассеиваются вдоль прямой (51), проходящей через начало координат. Ее угловой коэффициент  $\beta = 2,13537$ . Следовательно, согласно (47), параметр  $\beta = 1,0945$ . Параметр  $\alpha$ , вычисленный по формуле (43), равен  $\alpha = 0,03709$ .

Таким образом, выравнивающая кривая роста может быть записана в виде

$$y = 216696 e^{0,03709(x-1958)^{1,0945}}$$

Поскольку параметр  $\beta > 1$ , то темп роста числа названий книг и брошюр со временем увеличивается.

#### ЛИТЕРАТУРА

- Алексеев П.М. Частотный словарь английского подъязыка электроники // Статистика речи. - Л., 1968. - с. 161-161.
- Ешан Л.И. Опыт статистического описания научно-технического стиля румынского языка (на материале текстов по радиоэлектронике): Автореф. дис... канд. филол. наук. - Л., 1986. - 16 с.
- Калинин В.М. Некоторые статистические законы математической лингвистики // Проблемы кибернетики. Вып. II. - М., 1964 - с. 246-255.
- Калинина Е.А. Изучение лексико-статистических закономерностей на основе вероятностной модели // Статистика речи. - Л.: Наука, 1968. - с. 64-107.
- Кочеткова В.К., Скредина Л.М. Частотный словарь французского подъязыка электроники // Статистика речи. - Л., 1968. - с. 162-170.
- Михайлов А.И., Черный А.И., Гиляревский Р.С. Научные коммуникации и информатика. - М.: Наука, 1976. - 436 с.
- Нешиной В.В. Система непрерывных распределений как экономико-математическая модель // Проблемы макроэконометрического моделирования и прогнозирования: Тезисы докладов респ. конф. / АН Латв. ССР; ин-т эк-ки. - Рига: Зинатне, 1984. - с. 237-239.
- Нешиной В.В. Исследование ранговых распределений // НТИ. Сер.2. - 1985. - № 2. - с. 16-20.
- Нешиной В.В. Ранжирование слов по степени семантической нагрузки // НТИ. Сер.2. - 1986. - № 4. - с. 20-25.
- Пиотровский Р.Г., Алексеев П.М., Чернядьева Е.А. Статистика речи и закономерности языка // Тез. докл. межвуз. конф. на тему "Язык и речь" (27 ноября - 1 декабря) / МТИИ. - М., 1962. - с. 57-59.
- Частотный словарь русского языка. / Под ред. Л.Н. Засориной. - М.: Русский язык, 1977. - 936 с.
- Meier H. Deutsche Sprachstatistik. - Hildesheim, 1964.

MATHEMATICAL MODELS OF THE DICTIONARY  
EXPANSION AND DATA FLOWS

V.V. Neshitov

S u m m a r y

Four systems of the mathematical models for description of different growth curves are proposed, methods for parameter evaluation are given.

Indices of the analytical character of the language, word cohesion in a lexically uniform text and lexical similarity of two connected texts are introduced, these indices being independent of the text size.