

МОДЕЛИРОВАНИЕ
В ФИЗИКО-ТЕХНИЧЕСКИХ ИССЛЕДОВАНИЯХ

УДК 519.24

Е. Л. Кулешов, В. В. Крысанов, К. Какушо

(Владивосток, Россия – Кусацу, Киото, Япония)

РАСПРЕДЕЛЕНИЕ ВЕРОЯТНОСТЕЙ
ЧАСТОТЫ СЛОВА В ТЕКСТАХ*

Предлагается новая математическая модель формирования распределения вероятностей частоты слова в текстах, таких как английский, русский, английский гипертекст. Получено распределение вероятностей частот, обобщающее закон Парето. Представлен алгоритм вычисления параметров модели. Показано, что полученное распределение вероятностей обеспечивает существенно более высокую степень согласия с экспериментальными данными, чем распределение Парето.

Введение. Частота использования слова (число использований слова) в том или ином документе, тексте является важнейшей характеристикой, влияющей на динамику связи в современной Интернет-сети. Например, такие системы поиска информации, как Google, Yahoo, Altavista и т. п., используют алгоритмы, основанные на вычислении функции разрешающей способности (значимости) слова относительно текста, которая непосредственно определяется распределением вероятностей случайной величины [1–3]. Использование более точной модели формирования распределения вероятностей частоты позволяет повысить эффективность алгоритмов автоматической классификации текстов [4], а также алгоритмов поиска цифровых документов [5, 6].

Принято считать, что k – дискретная случайная величина с распределением вероятностей $P(k) = b/k^1$, где $b > 0$, $k = 1, 2, \dots$ и $P(k)$ – вероятность случайного события k . Это соотношение известно в литературе как распределение Парето [7], а также как второй закон Ципфа [8] – важнейший эмпирический закон, позволяющий моделировать сложные недетерминированные системы (например, телекоммуникационные сети [9]) и некоторые стохастические процессы (например, количество ролей, сыгранных актерами Голливуда в художественных фильмах [9], рост численности персонала компаний [10] и т. д.). В работах [11, 12] авторы выделяют три основ-

* Работа выполнена при финансовой поддержке Японского агентства науки и технологий (проект «Универсальный дизайн цифрового города»).

ных механизма, предположительно порождающих данное распределение вероятностей: 1) как выбор, основанный на предпочтениях (механизм впервые предложен в [13]), 2) как следствие оптимизационных процессов [14], а также 3) как следствие мультипликативных стохастических процессов различной природы [15]. В работе Саймона [13] показано, что частота слова в тексте, а также численность населения городов имеют распределение Парето. Аналогично Мандельброт в работе [14] дает теоретическое обоснование первого закона Ципфа: $q_n \sim dn^{-k}$, где n – ранг (порядковый номер) слова в последовательности слов, упорядоченных по убыванию частоты использования (наименьший ранг у слова, наиболее часто встречающегося в тексте); q_n – вероятность использования слова с рангом n ; d и k – константы. Известно, что существует взаимно однозначное преобразование между первым и вторым законами Ципфа: вывод первого из них влечет обоснование второго, и наоборот. Среди других исследований в этом направлении отметим математическую теорию эволюции Юла; цикл работ Чемпернауна [15] о закономерностях распределения доходов в экономике; подход, основанный на психофизиологических представлениях о деятельности головного мозга (А. Н. Лебедев) [16]; результат работы М. В. Арапова и др. в [16], полученный на основе точного учета целочисленности ранга слова.

Основной проблемой существующих аналитических моделей является наличие в них предположений, которые не всегда точно выполняются для исследуемых систем или процессов. Это несоответствие обычно указывается в качестве причины расхождений экспериментальных данных с результатами численного моделирования. Статистический анализ разнообразных данных (таких как частота слов в литературных произведениях на разных языках, индивидуальный и корпоративный доходы на примере различных государств, частота цитирования авторов и работ в научной литературе, численность персонала организаций и компаний, частота посещения различных Интернет-сайтов и Web-страниц и т. д. (см. обзор в [11, 12, 17])), выполненный многими авторами, вполне убедительно показывает, что закон Парето справедлив в основном в диапазоне больших k . Однако вместо поиска новых содержательных аналитических моделей исследования в данной области направлены, по сути, на подгонку результатов измерений (обычно путем предварительной фильтрации выборки) к одной из двух классических (Саймона или Мандельброта) моделей или же, наоборот, на подгонку этих моделей к данным путем ввода дополнительных параметров или компонент, которые не имеют какого-либо физического смысла [18]. Адекватная математическая модель должна, по всей видимости, учитывать, во-первых, дискретность случайной величины n и, во-вторых, психофизиологический фактор формирования естественного языка.

В данной работе предлагается новая математическая модель механизма формирования распределения вероятностей частоты слова в текстах, таких как английский, русский, а также английский гипертекст (актуальная задача моделирования связи в Интернет-сети, которая не имеет решения в рамках классических подходов [19]). Получено распределение вероятностей частот, обобщающее закон Парето. Представлен алгоритм вычисления параметров модели. На примере разнообразных данных (смесь английских текстов, смесь гипертекстов, смесь текстов на русском языке) показано, что полученное распределение вероятностей обеспечивает существенно более высокую степень согласия с экспериментальными данными, чем распределение Парето.

1. Модель механизма формирования распределения вероятностей частоты слова. Для обозначения некоторого концепта (понятия) в естественных языках используется не единственное слово, а, как правило, множество слов (например, синонимы или слова, заимствованные из другого языка). Пусть p_k – вероятность использовать слово k раз для обозначения фиксированного концепта, т. е. p_k – условная вероятность события k при условии, что рассматривается только фиксированный концепт. Будем полагать, что среди всех распределений p_k с фиксированным математическим ожиданием $a = \sum_k k p_k$ реализуется такое распределение, для которого энтропия $H = - \sum_k p_k \ln p_k$ достигает своего максимального значения. Поиск максимума функции H с условиями $\sum_k p_k = 1$ и $\sum_k k p_k = a$ сводится к задаче поиска безусловного максимума функции Лагранжа

$$- \sum_k p_k \ln p_k - \lambda \left(\sum_k p_k - 1 \right) - \mu \left(\sum_k k p_k - a \right) \quad (1)$$

с множителями λ и μ , решение которой находится из системы уравнений: $\ln p_k = -1 - \lambda - \mu k$, $\sum_k p_k = 1$, $\sum_k k p_k = a$ [20]. Из первого и второго уравнений следует $p_k = P(k) = \frac{1}{Z} e^{-\lambda - \mu k}$, где

$$P(k) = \frac{1}{Z} (e^{-\mu})^k, \quad k = 1, 2, \dots, \quad (2)$$

– условная вероятность события k при фиксированном λ . Таким образом, параметр μ является статистической характеристикой рассматриваемого концепта. Из (2) следует

$$a = \sum_{k=1}^{\infty} k p_k = e^{-\lambda} / (e^{-\mu} - 1) = 1/p_1.$$

Математическое ожидание a числа использований слова существенно больше единицы, поэтому параметры p_1 и μ существенно меньше единицы. При малом a получаем $\mu \approx 1/a$. Величина, обратная среднему числу использования слова, пропорциональна μ – среднему времени использования слова носителем языка для обозначения рассматриваемого концепта. Отметим, что экспоненциальное распределение частоты можно получить другим способом [21].

Параметр μ будем рассматривать как непрерывную случайную величину с плотностью $f(\mu)$, следовательно, безусловное распределение вероятностей частоты использования слова имеет вид

$$P(k) = \int_0^{\infty} f(\mu) P(k|\mu) d\mu, \quad k = 1, 2, \dots \quad (3)$$

Для естественного языка измеряемая величина является смесью нескольких случайных величин, которые отражают влияние подсистем этого языка, например подсистемы служебных слов и подсистемы смысловых слов [22,

23]. При этом плотность распределения $f(k) = \sum_{i=1}^N c_i f_i(k)$, где N – число компонент распределения (подсистем); c_i – вероятность обнаружить i -ю компоненту в наблюдаемой смеси; $f_i(k)$ – плотность вероятности случайной величины k в i -й компоненте. Тогда

$$P(k) = \sum_{i=1}^N c_i P_i(k), \quad (4)$$

$$P_i(k) = \int_0^{\infty} f_i(k/d) d^{\alpha} d, \quad k = 1, 2, \dots, \quad (5)$$

где $P_i(k)$ – распределение вероятности частоты использования слова в i -й компоненте.

Каждая из функций $f_i(k)$ – это, по сути, плотность вероятности величины k – среднего времени использования слова для обозначения фиксированного концепта. Результаты многих физиологических экспериментов показывают, что k линейно зависит от среднего времени реакции, которое является гамма-распределенной случайной величиной [24]. Поэтому полагаем $f_i(k) = \frac{b_i^{\alpha} k^{\alpha-1} e^{-b_i k}}{\Gamma(\alpha)}$, $\alpha > 0$, где $\Gamma(\alpha)$ – гамма-функция, а $b_i > 0$ и $\alpha > 0$ – параметры распределения. Подстановка этого распределения и формулы (2) в (5) приводит к выражению

$$P_i(k) = \frac{b_i^{\alpha}}{(k+1) b_i)^{\alpha}} = \frac{b_i^{\alpha}}{(k+b_i)^{\alpha}}, \quad k = 1, 2, \dots \quad (6)$$

Таким образом, распределение вероятностей частоты использования слова определяется соотношениями (4), (6), которые получены на основе следующих положений предлагаемой математической модели механизма формирования частоты.

1. Для фиксированного концепта реализуется распределение вероятностей частоты использования слова с максимальным значением энтропии.

2. Среднее время использования слова для обозначения концепта зависит от среднего времени реакции и есть гамма-распределенная случайная величина.

3. Естественный язык как система передачи информации может состоять из нескольких статистически различных подсистем.

2. Асимптотика распределения частоты. Если предположить, что соотношения (4), (6) хорошо согласуются с экспериментальными данными для всех значений $k = 1, 2, \dots$, тогда в асимптотике при больших k эти соотношения должны переходить в распределение Парето, поскольку, как отмечалось, закон Парето хорошо согласуется с экспериментальными данными для больших k . Покажем, что действительно из соотношений (4), (6) при некоторых условиях следует распределение Парето. Пусть число компонент $N = 1$, тогда $c_1 = 1$ и сумма (4) содержит единственное слагаемое $P(k) = P_1(k)$. При этом в формуле (6) у параметров b_i, α опустим индекс i , поскольку рассматривается единственное значение $i = 1$. Пусть выполняется условие $k \gg b$, тогда $z = 1/(k+b) \approx 1/k$. Разлагая в ряд Тейлора по малому параметру

z функцию $(1 - z)^{-b}$, имеем $(1 - z)^{-b} = 1 + bz + \dots$. С учетом этого первое слагаемое (6) можно преобразовать следующим образом:

$$\frac{b}{(k-1-b)} = \frac{b}{(k-b)(1-z)} = \frac{b}{(k-b)} \cdot \frac{1}{1-z} \quad (7)$$

Подставим (7) в (6), тогда

$$P_1(k) = \frac{b}{(k-b)^{b+1}}, \quad k = b+1. \quad (8)$$

Пусть выполняется дополнительное условие $k = b$, тогда из (8) следует распределение Парето $P_1(k) = b/k^{b+1}$. Отметим, что при $b=1$ первому условию $k = b+1$ удовлетворяют значения $k=1$, а при $b=1$ – значения $k=1, 2, \dots$. Второму условию $k = b$ при любом b удовлетворяют только значения $k=1$. Поэтому второе условие более сильное, чем первое. Распределение Парето можно получить непосредственно из (6) при условии $k = b$. Для этого достаточно разложить в ряд Тейлора первое слагаемое (6) по малому параметру $(b-1)/k$ и второе слагаемое по малому параметру b/k .

Соотношение (8) допускает следующую интерпретацию. Пусть рассматривается как непрерывная случайная величина, тогда аналогом формулы (2) будет экспоненциальная плотность распределения вероятностей e^{-k} , $k \geq 0$. При этом вычисление интеграла (5) приводит к соотношению (8). Таким образом, (8) можно рассматривать как плотность распределения вероятностей непрерывной случайной величины для $k \geq 0$. Такая плотность имеет максимальное значение $P_1(0) = 1/b$, является монотонной убывающей

функцией и удовлетворяет условию нормировки $\int_0^{\infty} P_1(k) dk = 1$. Пусть теперь

имеются две плотности $P_1(k)$ и $P_2(k)$ вида (8) с параметрами b_1, b_1 и b_2, b_2 соответственно. Тогда условие $b_1/b_1 = b_2/b_2$ эквивалентно $P_1(0) = P_2(0)$, откуда следует, что при больших k должно выполняться обратное соотношение $P_1(k) < P_2(k)$, поскольку в противном случае нарушается условие нормировки. Это свойство функции (8) полезно для последующей интерпретации результатов анализа экспериментальных данных.

3. Оценки максимального правдоподобия. Проверка соотношений (4), (6) на соответствие экспериментальным данным заключается в оценивании параметров $N, b_1, \dots, b_N, \mu_1, \dots, \mu_N$ и последующем использовании критерия согласия теоретического распределения $P(k)$ (с параметрами, равными их оценкам) и эмпирического распределения $W(k)$ частоты использования слова. Однако первая задача оценивания параметров распределения (4) в такой общей постановке очень сложна и вряд ли может быть решена известными методами теории оценивания [25]. Так, метод моментов здесь вообще неприменим, поскольку для распределения (6) существует небольшое число моментов, недостаточное для построения системы уравнений относительно оценок. Использование критерия χ^2 осложняется проблемой минимизации нелинейной функции $2N-1$ переменных. Решение системы уравнений максимального правдоподобия становится проблематичным даже при $N=1$. Эти

трудности приводят к необходимости поиска приближенных решений задачи оценивания.

Рассмотрим простой вариант приближенного решения в предположении, что распределение вероятностей частоты слова для $k = 1, 2, \dots$ задается соотношением (8). Имеются следующие доводы в пользу такого варианта.

1. Если верны соотношения (4), (6), тогда соотношение (8) верно для более широкого диапазона значений k , чем распределение Парето, так как последнее следует из (6) при более сильных ограничениях, чем (8). Поэтому можно ожидать, что соотношение (8) точнее описывает распределение частоты слова, чем распределение Парето.

2. Соотношение (8) имеет простые аналитические свойства, что позволяет получить практический алгоритм вычисления оценок максимального правдоподобия.

Пусть x_1, \dots, x_n – выборка объема n , полученная при выполнении n измерений случайной величины, тогда функция правдоподобия для распределения (8) имеет вид

$$P_1(x_i) = \frac{b}{(x_i + b)^{b+1}}. \quad (9)$$

Введем обозначение

$$L = \ln \prod_{i=1}^n P_1(x_i) = - \sum_{i=1}^n [\ln b + (x_i + b)^{-1}] \quad (10)$$

и составим систему уравнений максимального правдоподобия: $L'_{b'} = 0$, $L'_{b} = 0$, решение которой будет определять оценки параметров b' , b . Пусть

$$t_1 = \frac{1}{n} \sum_{i=1}^n \frac{b}{x_i + b}; \quad t_2 = \frac{1}{n} \sum_{i=1}^n \ln \frac{b}{x_i + b}, \quad (11)$$

тогда система уравнений максимального правдоподобия приводится к виду

$$\frac{1}{b'} = t_2; \quad \frac{t_1}{1 - t_1} = b'. \quad (12)$$

Исключив параметр b' из соотношений (12), получим

$$1 - t_1 = t_1 t_2 \quad (13)$$

– уравнение относительно оценки параметра b , которое несложно решить численным методом. Подстановка этого решения в одно из соотношений (12) позволяет вычислить оценку параметра b' . Таким образом, выражения (12), (13) определяют решение системы уравнений максимального правдоподобия.

Для проверки соотношения (8) на соответствие экспериментальным данным была сформирована последовательность частоты использования слова в статьях политической тематики на английском языке, выбранных на сайте издательства “Associated Press”. При подсчете частоты слова игнорировались различия между заглавным и строчным регистрами в написании букв.

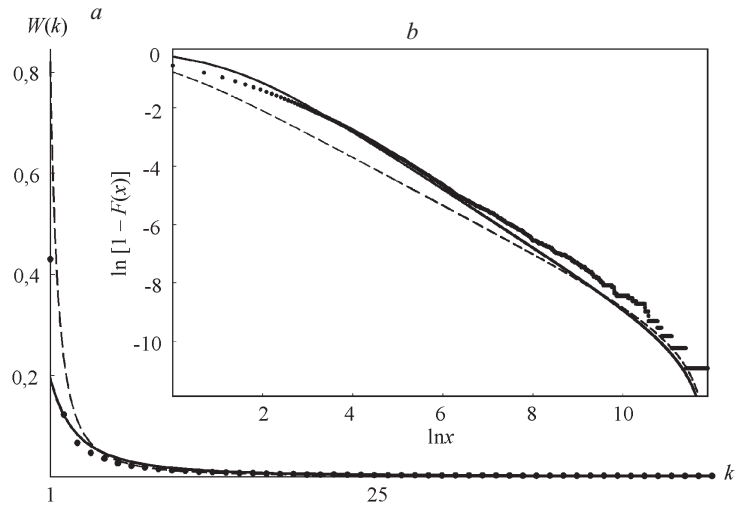


Рис. 1

Во всех экспериментах какой-либо иной предварительной обработки текстов, в том числе для учета морфологической изменчивости слов, не проводилось. Общий объем обработанных текстов составил 7,3 Мбайт при количестве слов 10^6 . На рис. 1, а точками представлено эмпирическое распределение $W(k)$ частоты использования слова, пунктиром – распределение Парето и сплошной линией – распределение вида (8). Оценки максимального правдоподобия параметров распределения Парето были получены с использованием программы BestFit 4.5 “Palisade Corporation” для научных и инженерных статистических расчетов. На рис. 1, b в двойном логарифмическом масштабе представлены соответствующие интегральные функции распределения вероятностей для целого аргумента, например эмпирическая функция распределения $F(k) = \sum_{i=1}^k W(i)$, $k = 1, 2, \dots$. Необходимо отметить,

что наблюдаемое расхождение классической модели с эмпирическими данными хорошо известно в литературе [11] и не зависит от параметров текста, таких как объем, язык, однородность (в смысле, что это текст одного автора или подборка текстов разных авторов) и т. д.

4. Оценки параметров второй компоненты. Результаты анализа экспериментальных данных показывают, что соотношение (8) точнее описывает распределение вероятностей частоты слова, чем закон Парето, однако также не обеспечивает полного соответствия с эмпирическим распределением, особенно в области малых k . Видимо, это обусловлено тем, что естественный язык представляет собой смесь, как минимум, двух статистически различных подсистем: служебных и смысловых слов [22]. Поэтому можно ожидать существенного повышения степени согласия эмпирического $W(k)$ и теоретического $P(k)$ распределений, если использовать для последнего выражение (4) при $N = 2$. Статистическая различимость подсистем при этом означает существенное различие распределений $P_1(k)$ и $P_2(k)$ в формуле (4), что приводит к условию $P_1(k) \neq P_2(k)$ в области малых k и $P_1(k) \approx P_2(k)$ для больших k . Таким образом, при больших k из (4) следует $P(k) \approx c_2 P_2(k)$.

Это позволяет сформулировать критерий выбора оценок параметров c_2, σ_2, b_2 в виде условия

$$\min_{c_2, \sigma_2, b_2} \sum_{k=n_1}^{n_2} [\ln W(k) - \ln c_2 P_2(k)]^2. \quad (14)$$

Здесь число n_1 выбирается столь большим, чтобы при $k = n_1$ выполнялось неравенство $P_1(k) > P_2(k)$ и в соотношении (4) первым слагаемым можно было пренебречь. Второе число n_2 зависит от объема выборки и выбирается не слишком большим, так чтобы при $k = n_2$ эмпирическое распределение $W(k)$ существенно отличалось от нуля. При этом для $P_2(k)$ справедлива асимптотика (8). Подставим (8) в (14) и для упрощения обозначений в последующих преобразованиях опустим индексы у параметров c_2, σ_2, b_2 , тогда

$$\min_{c, \sigma, b} \sum_{k=n_1}^{n_2} [\ln W(k) - \ln c - \ln \sigma - \ln b - (1 - \sigma) \ln(k - b)]^2. \quad (15)$$

Таким образом, задача оценивания параметров второй компоненты сводится к поиску минимума функции (15) трех аргументов c, σ, b . Значения аргументов в точке минимума принимаются за оценки параметров второй компоненты. Необходимые условия минимума функции нескольких аргументов в данном случае сводятся к системе трех уравнений [20]:

$$\frac{\partial}{\partial c} = 0; \quad \frac{\partial}{\partial \sigma} = 0; \quad \frac{\partial}{\partial b} = 0. \quad (16)$$

Пусть $n = n_2 - n_1 + 1$ – число слагаемых в сумме (15) и

$$z = \ln c - \ln \sigma - \ln b. \quad (17)$$

Тогда первое уравнение системы (16) принимает вид

$$\frac{1}{n} \sum_{k=n_1}^{n_2} \ln W(k) - z - (1 - \sigma) \frac{1}{n} \sum_{k=n_1}^{n_2} \ln(k - b) = 0. \quad (18)$$

С учетом этого равенства второе уравнение системы (16) преобразуется к выражению

$$\frac{1}{n} \sum_{k=n_1}^{n_2} \ln W(k) \ln(k - b) - z \frac{1}{n} \sum_{k=n_1}^{n_2} \ln(k - b) - (1 - \sigma) \frac{1}{n} \sum_{k=n_1}^{n_2} \ln^2(k - b) = 0. \quad (19)$$

Третье уравнение (16) также с учетом равенства (18) сводится к выражению

$$\frac{1}{n} \sum_{k=n_1}^{n_2} \frac{\ln W(k)}{k - b} - z \frac{1}{n} \sum_{k=n_1}^{n_2} \frac{1}{k - b} - (1 - \sigma) \frac{1}{n} \sum_{k=n_1}^{n_2} \frac{\ln(k - b)}{k - b} = 0. \quad (20)$$

Соотношения (18)–(20) образуют систему трех уравнений относительно неизвестных z , b , решение которой позволяет найти параметр c из равенства (17).

Введем обозначения:

$$s_1 = \frac{1}{n} \sum_{k=1}^{n_2} \frac{1}{k-b}; \quad s_2 = \frac{1}{n} \sum_{k=1}^{n_2} \frac{\ln(k-b)}{k-b}; \quad s_3 = \frac{1}{n} \sum_{k=1}^{n_2} \frac{\ln W(k)}{k-b}, \quad (21)$$

$$z_1 = \frac{1}{n} \sum_{k=1}^{n_2} \ln(k-b); \quad z_2 = \frac{1}{n} \sum_{k=1}^{n_2} \ln^2(k-b), \quad (22)$$

$$z_3 = \frac{1}{n} \sum_{k=1}^{n_2} \ln W(k) \ln(k-b); \quad z_4 = \frac{1}{n} \sum_{k=1}^{n_2} \ln W(k). \quad (23)$$

Выразим z из уравнения (18) и подставим его в (19) и (20), тогда в принятых обозначениях система трех уравнений имеет следующий вид:

$$z = z_4 + (1 - c)z_1, \quad (24)$$

$$z_3 = z_1 z_4 + (1 - c)(z_2 - z_1^2) = 0, \quad (25)$$

$$s_3 = s_1 z_4 + (1 - c)(s_2 - s_1 z_1) = 0. \quad (26)$$

Из (25) находим

$$1 = \frac{z_3 - z_1 z_4}{z_2 - z_1^2}. \quad (27)$$

В соответствии с формулами (21)–(23) величины $s_1, s_2, s_3, z_1, z_2, z_3$ являются функциями аргумента b и z_4 – число. Поэтому соотношение (27) определяет параметр c как функцию аргумента b . Подставим (27) в (26) и получим уравнение

$$s_3 - s_1 z_4 + (s_2 - s_1 z_1)(z_3 - z_1 z_4)/(z_2 - z_1^2) = 0 \quad (28)$$

относительно неизвестного b . Его решение является оценкой параметра b_2 . При $b = b_2$ соотношение (27) определяет оценку параметра c_2 , подстановка $c_2, b = b_2$ в (24) – величину z , и затем из соотношения (17) вычисляется оценка c_2 .

5. Оценки параметров первой компоненты. Подставим полученные параметры c_2, b_2 в формулу (6) и вычислим $P_2(k)$ для $k = 1, 2, \dots$. Затем найдем

$$W_1(k) = \frac{W(k) - c_2 P_2(k)}{1 - c_2}, \quad k = 1, 2, \dots, \quad (29)$$

– эмпирическое распределение частоты слова первой компоненты. Здесь $1 - c_2 - c_1$ – оценка вероятности того, что слово принадлежит первой компоненте. Приравнявая теоретическое (6) и эмпирическое распределения

$P_1(k) = W_1(k)$ для $k = 1, 2$, получаем систему двух уравнений относительно неизвестных c_1, b . Из первого уравнения ($k = 1$) следует

$$\frac{\ln[1 - W_1(1)]}{\ln b - \ln(b - 1)}. \quad (30)$$

Подставим (30) во второе уравнение ($k = 2$) и получим

$$\frac{b}{b - 1} \frac{\ln[1 - W_1(1)]}{\ln b - \ln(b - 1)} = \frac{b}{b - 2} \frac{\ln[1 - W_1(1)]}{\ln b - \ln(b - 1)} W_1(2) \quad (31)$$

– уравнение относительно b . Его решение принимаем за оценку параметра b_1 . Подстановка $b = b_1$ в соотношение (30) определяет оценку параметра c_1 .

Представленный алгоритм использовался для вычисления параметров теоретического распределения вида (4), (6) по той же выборке частот, которая рассматривалась в разд. 3. На рис. 2, а показаны эмпирическое распределение $W(k)$ (точки) и теоретическое распределение $P(k)$ (сплошная линия), а на рис. 2, б – эмпирическая функция распределения $F(k)$ (точки) и теоретическая функция распределения для целого аргумента (сплошная линия).

Такие же результаты получены для английского гипертекста (выборка Web-страниц сайта www.cbsnews.com, сформированная случайным образом; объем 11,4 Мбайт, число слов $1,8 \cdot 10^6$) и текста на русском языке (подборка коротких литературных произведений разных авторов, произвольным образом выбранных на сайте www.lib.ru; объем 3,1 Мбайт, число слов $0,45 \cdot 10^6$). Основные параметры $c_1, c_2, c_1/b_1, c_2/b_2$ распределения частоты для английского текста равны 0,525, 0,475, 0,805, 0,136; для гипертекста – 0,549, 0,451, 0,572, 0,123; для русского текста – 0,673, 0,327, 0,978, 0,363 соответственно. Отметим, что для английского текста и гипертекста расчетные значения вероятностей c_1 и c_2 хорошо согласуются с известной

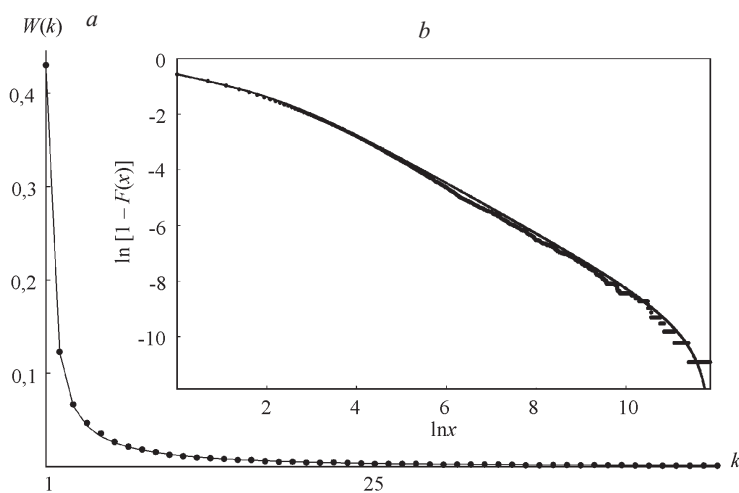


Рис. 2

статистикой смысловых и служебных слов, на которые в современном английском языке приходится соответственно 50–55 % и 45–50 % объема текста [22, 11]. Этот факт представляется весомым аргументом в пользу предлагаемой модели.

Отметим, что критерий оптимальности (14) – это частный случай метода наименьших квадратов [26, 27], где функция $\ln c_2 P_2(k)$ дифференцируема по всем параметрам. Для большого объема выборки случайные остатки $\ln W(k) - \ln c_2 P_2(k)$ имеют нулевые средние и одинаковые дисперсии. Поэтому условие (14) определяет состоятельные оценки. Оценки параметров первой компоненты (29)–(31) также состоятельны, поскольку вычисляются по методу моментов: приравниванием эмпирического и теоретического распределений вероятностей для $k = 1, 2$.

Заключение. В работе получено распределение вероятностей частоты использования слова в текстах, обобщающее закон Парето, а также разработан алгоритм вычисления его параметров. На примере разнообразных данных показано, что полученное распределение вероятностей обеспечивает существенно более высокую степень согласия с экспериментальными данными, чем распределение Парето. Этот результат представляется принципиально важным с точки зрения создания разнообразных сетевых мультимедийных приложений [19], а также систем поиска информации в Интернет-сети [28], так как полученная аналитическая модель позволяет с хорошей точностью вычислить важнейшие характеристики процесса передачи информации. Можно ожидать, что полученное распределение будет полезно для анализа работы телекоммуникационных сетей [9, 29, 30] и различных социальных, в том числе экономических, систем [10, 31].

СПИСОК ЛИТЕРАТУРЫ

1. Van Rijsbergen C. J. Information Retrieval. London: Butterworths, 1979.
2. Losee R. M. Term dependence: A basis for Luhn and Zipf models // Journ. Amer. Soc. Inform. Sci. and Technol. 2001. 52, N 12. P. 1019.
3. Brin S., Page L. The anatomy of a large-scale hypertextual Web search engine // Comput. Networks and ISDN Systems. 1998. 30, N 1–7. P. 107.
4. Peters C., Koster C. Uncertainty and term selection in text categorization // Intern. Journ. Uncertainty, Fuzziness and Knowledge-Based Systems. 2003. 11, N 1. P. 115.
5. Amati G., Van Rijsbergen C. J. Probabilistic models of information retrieval based on measuring the divergence from randomness // ACM Trans. on Inform. Systems (TOIS). 2002. 20, N 4. P. 357.
6. Amati G., Van Rijsbergen C. J. Term frequency normalization via Pareto distributions // Springer Lecture Notes in Comput. Sci. 2002. 2291. P. 183.
7. Pareto V. Cours d'Economie Politique. Lausanne – Paris: Rouge, 1897.
8. Zipf G. K. Human Behavior and the Principle of Least Effort. Cambridge, Mass.: Addison-Wesley, 1949.
9. Barabasi A., Albert R. Emergence of scaling in random networks // Sci. Magazine. 1999. 286. P. 509.
10. Axtell R. L. Zipf distribution of U.S. firm sizes // Sci. Magazine. 2001. 293. P. 1818.
11. <http://www.ram-verlag.de> (Glottometrics, Special issue to honor G. K. Zipf. 2002. N 4.)

12. Mitzenmacher M. A brief history of generative models for power law and lognormal distributions // *Internet Mathem.* 2003. 1, N 2. P. 226.
13. Simon H. A. On a class of skew distribution functions // *Biometrika.* 1955. N 42. P. 425.
14. Mandelbrot B. B. An informational theory of the statistical structure of languages // *Commun. Theory.* Betterworth. 1953. P. 486.
15. Champernowne D. *The Distribution of Income.* Cambridge: Cambridge University Press, 1973.
16. Тулдава Ю. А. Проблемы и методы квантитативно-системного исследования лексики. Таллин: Валгус, 1987.
17. Shide N., Batty M. Power Law Distributions in Real and Virtual Worlds // <http://www.isoc.org/inet2000> (Inet 2000 Proc., Internet Soc. 2000.)
- 18.
- 19.
20. Бронштейн И. Н., Семендяев К. А. *Справочник по математике.* М.: Наука, 1986.
21. Kryssanov V. V., Kakusho K., Kuleshov E. L., Minoh M. Modeling hypermedia-based communication // *Intern. Journ. Inform. Sci.* 2004; <http://authors.elsevier.com/sd/article/S0020025504002385>
22. Balasubrahmanyam V. K., Naranan S. Quantitative linguistics and complex system studies // *Journ. Quant. Linguistics.* 1996. N 3. P. 177.
23. Naranan S., Balasubrahmanyam V. K. Information theoretic models in statistical linguistics. Pt. I, II // *Current Sci.* 1992. N 63. P. 261, 297.
24. Luce R. D. *Response Times. Their Role in Inferring Elementary Mental Organization.* New York: Oxford University Press, 1986.
25. Королук В. С., Портенко Н. И., Скороход А. В., Турбин А. Ф. *Справочник по теории вероятностей и математической статистике.* М.: Наука, 1985.
26. Рао С. Р. *Линейные статистические методы и их применения.* М.: Наука, 1968.
27. Айвазян С. А., Енюков И. С., Мешалкин Л. Д. *Прикладная статистика. Основы моделирования и первичная обработка данных.* М.: Финансы и статистика, 1983.
28. Bates M. Indexing and access for digital libraries and the Internet: Human, database, and domain factors // *Journ. Amer. Soc. Inform. Sci.* 1998. 49, N 13. P. 1185.
29. Huberman B. A., Pirolli P. L. T., Pitkow J. E., Lukose R. M. Strong regularities in World Wide Web surfing // *Sci. Magazine.* 1998. N 280. P. 95.
30. Adamic L. A., Huberman B. A. Power-law distribution of the World Wide Web // *Sci. Magazine.* 2000. 287. P. 2115a.
31. Stanley M. H. R., Amaral L. A. N., Buldyrev S. V. et al. Scaling behavior in the growth of companies // *Nature.* 1996. 379. P. 804.

Дальневосточный государственный университет,
 Рицумеикан университет, Киото университет,
 E-mail: kuleshov@lemoi.phys.dvgu.ru

Поступила в редакцию
 8 июня 2004 г.