

М. А. Лифшиц, И. М. Лялинов

ДВЕ ПРЕДЕЛЬНЫЕ ТЕОРЕМЫ О ПЕРЕСЕЧЕНИЯХ СЛУЧАЙНЫХ МНОЖЕСТВ ЦИПФА

§1. ВВЕДЕНИЕ И ОСНОВНЫЕ РЕЗУЛЬТАТЫ

1.1. Случайные множества Ципфа. Эмпирический закон Ципфа (G. K. Zipf [7,8], см. также [5]) заключается в том, что самое употребительное слово языка встречается вдвое чаще следующего по частоте, втрое чаще следующего, и т.д. Ципф также применял свой закон к таким объектам экономики, как население городов, размеры индивидуальных состояний и др.¹

Среди относительно недавних работ упомянем приложения в лингвистике [9, 12] и ряд работ, посвящённых исследованию вероятностных и комбинаторных механизмов, приводящих к закону Ципфа [13–17].

Пусть имеется счётное множество объектов $\mathcal{W} = \{w_1, w_2, \dots\}$. Случайным множеством Ципфа с параметрами $a \in (0, 1]$, $\alpha \in (0, 1]$ назовём случайное подмножество $A \subset \mathcal{W}$, для которого события $\{w_m \in A\}$ независимы и $\mathbb{P}(w_m \in A) = \frac{a}{m^\alpha}$.

В дальнейшем для краткости обозначений мы отождествляем \mathcal{W} с множеством натуральных чисел \mathbb{N} , а w_m с m .

Отметим, что условие $a \leq 1$ совершенно несущественно и выбрано только для удобства работы с произвольными натуральными индексами. Все результаты статьи остаются верными для случайных множеств, удовлетворяющих асимптотическому условию

$$\lim_{m \rightarrow \infty} m^\alpha \mathbb{P}(w_m \in A) = a$$

с произвольным $a > 0$.

Ключевые слова: предельные теоремы, максимальные пересечения, множества Ципфа, случайные множества, наиболее редкий элемент пересечения.

Работа поддержана грантом РФФ No. 21-11-00047.

¹В лингвистическом аспекте Ципф существенно опирался на данные, накопленные стенографами, в числе которых Ж.-Б. Эсту (J.-B. Estoup [2]) уже в начале XX века отмечал гиперболический характер наблюдаемых частот, см. [4, 6], а применительно к размеру городов закон сформулировал ещё физик Ф. Ауэрбах (F. Auerbach, 1913, [1]).

1.2. Максимальные элементы пересечений множеств Ципфа.

Пусть A – случайное множество Ципфа с параметрами $a \in (0, 1]$, $\alpha \in (0, 1]$. Очевидно, что такое множество почти наверное бесконечно (но его пересечение с независимым случайным множеством Ципфа с параметрами $x \in (0, 1]$, $\beta \in (1 - \alpha, 1]$ почти наверное конечно). Поэтому, прежде всего, нас будут интересовать характеристики пересечений случайных множеств Ципфа с указанным выше ограничением на параметры.

Будем писать $c_n \sim d_n$, если $\lim_{n \rightarrow \infty} \frac{c_n}{d_n} = 1$ и $c_n \preceq d_n$, если выполнено $\limsup_{n \rightarrow \infty} \frac{c_n}{d_n} \leq 1$.

Задача о максимальных характеристиках пересечений множеств возникает в математическом моделировании функционирования баз данных. Представим себе, что база данных состоит из документов X_1, \dots, X_n , каждый из которых можно охарактеризовать набором слов (или ключевых слов), и имеется запрос A , который также характеризуется набором слов. Релевантность документа X запросу A можно определять различными способами. В работах [3, 11] релевантность понимается как количество общих слов в запросе и документе. В них получены законы больших чисел для максимальной меры релевантности при объёме базы, стремящемся к бесконечности.

В данной работе релевантность определяется другим естественным способом – как номер самого редкого (имеющего максимальный номер) общего слова в A и X . Задачей выполнения запроса является нахождение наиболее релевантного документа, т.е.

$$\operatorname{argmax}_{1 \leq j \leq n} V_j, \quad (1)$$

где через V_j обозначен максимальный элемент множества $A \cap X_j$.

Разумеется, представляет интерес и величина самого редкого общего элемента

$$V^n := \max_{1 \leq j \leq n} V_j,$$

например, для того, чтобы оценивать качество приближённых решений задачи (1). Отметим, что задача о максимальном пересечении имеет многочисленные приложения: в интернет-поиске, в рекомендательных системах, онлайн-рекламе и т.д.

Следующая предельная теорема характеризует распределение самого редкого общего элемента пересечения в модели, где и документы, и запрос рассматриваются как случайные множества Ципфа.

Теорема 1. Пусть $a, x \in (0, 1]$, $\alpha \in (0, 1]$, $\beta \in (1 - \alpha, 1]$.

Тогда при $\alpha < 1$ для любого $r \geq 0$ и $n \rightarrow \infty$ выполнено

$$\mathbb{P}\left(V^n < r n^{\frac{1}{\alpha+\beta-1}}\right) \rightarrow F_{\alpha,\beta}(r) := \exp\left(-\frac{ax}{(\alpha+\beta-1)r^{\alpha+\beta-1}}\right),$$

а при $\alpha = 1$ для любого $r \geq 0$ и $n \rightarrow \infty$ выполнено

$$\mathbb{P}\left(V^n < r n^{\frac{1}{\beta}}\right) \rightarrow F_{1,\beta}(r) := \exp\left(-\frac{aI(x/r^\beta)}{\beta}\right),$$

где

$$I(v) := \int_0^v \frac{1 - e^{-w}}{w} dw, \quad v \geq 0.$$

Отметим, что при $\alpha < 1$ предельное распределение $F_{\alpha,\beta}$ есть не что иное, как распределение Фреше, хорошо известное в теории предельных теорем для максимумов независимых одинаково распределённых величин [10]. Это не слишком удивительно, так как величины V_j при фиксированном A независимы и одинаково распределены. Однако распределение $F_{1,\beta}$ совсем не относится к классическим. Авторы сталкиваются с ним впервые.

Представляют интерес и меры релевантности, промежуточные между номером самого редкого общего элемента и количеством общих элементов. Пусть $(q_m)_{m \geq 1}$ – последовательность неотрицательных чисел (весов). Определим связанную с ними интегральную меру релевантности документа X запросу A как

$$Q(A, X) := \sum_{m \in A \cap X} q_m.$$

Если $q_m = 1$ при всех $m \geq 1$, то $Q(A, X)$ – количество общих элементов. Если же $q_m \nearrow \infty$ достаточно быстро, то асимптотическое поведение максимальной интегральной меры релевантности определяется асимптотикой самого редкого общего элемента. Сейчас мы установим этот факт для экспоненциально растущих весов.

Пусть $b > 0$, $q_m := e^{bm}$, $Q_j := Q(A, X_j)$, $Q^n := \max_{1 \leq j \leq n} Q_j$.

Теорема 2. Пусть $a, x \in (0, 1]$, $\alpha \in (0, 1]$, $\beta \in (1 - \alpha, 1]$.

Тогда для любого $r \geq 0$ при $n \rightarrow \infty$ выполнено

$$\mathbb{P}(\log Q^n \leq br n^\gamma) \rightarrow F_{\alpha,\beta}(r),$$

где функции распределения $F_{\alpha,\beta}(\cdot)$ определены в теореме 1 и

$$\gamma := (\alpha + \beta - 1)^{-1}.$$

В заключение отметим, что случай степенных весов ещё предстоит исследовать.

Доказательство теоремы 1. Пользуясь независимостью A и X_j , а также независимостью появления отдельных элементов в множествах Цифа, преобразуем интересующую нас вероятность в произведение вероятностей: для любого N верно

$$\begin{aligned} \mathbb{P}(V^n < N) &= \mathbb{P}\left(m \notin \bigcup_{j=1}^n (A \cap X_j), \forall m \geq N\right) \\ &= \prod_{m \geq N} \mathbb{P}\left(m \notin \bigcup_{j=1}^n (A \cap X_j)\right) \\ &= \prod_{m \geq N} \left(1 - \mathbb{P}\left(m \in \bigcup_{j=1}^n (A \cap X_j)\right)\right) \\ &= \prod_{m \geq N} \left(1 - \mathbb{P}(m \in A) \mathbb{P}\left(m \in \bigcup_{j=1}^n X_j\right)\right) \\ &= \prod_{m \geq N} \left(1 - \mathbb{P}(m \in A) \left(1 - \mathbb{P}\left(m \notin \bigcup_{j=1}^n X_j\right)\right)\right) \\ &= \prod_{m \geq N} \left(1 - \frac{a}{m^\alpha} \left(1 - \left(1 - \frac{x}{m^\beta}\right)^n\right)\right). \end{aligned} \quad (2)$$

Будем оценивать это произведение.

Пусть сначала $\alpha < 1$. Фиксируем $r > 0$ и положим $N := r n^{\frac{1}{\alpha+\beta-1}}$.

При $\alpha < 1$ в зоне $\{m : m \geq N\}$ верно $m^\beta \gg n$. Поэтому

$$1 - \left(1 - \frac{x}{m^\beta}\right)^n \sim \frac{xn}{m^\beta}$$

и

$$\begin{aligned} \log \mathbb{P}(V^n < N) &\sim - \sum_{m \geq N} \frac{anx}{m^{\alpha+\beta}} \sim - \frac{anx}{(\alpha + \beta - 1)N^{\alpha+\beta-1}} \\ &= - \frac{ax}{(\alpha + \beta - 1)r^{\alpha+\beta-1}}. \end{aligned}$$

Последнее соотношение эквивалентно утверждению теоремы при $\alpha < 1$.

В случае $\alpha = 1$ для фиксированного $r > 0$ положим $N := rn^{\frac{1}{\beta}}$.

Нам нужно показать, что величина $p_n := -\log \mathbb{P}(V^n < N)$ имеет соответствующий предел. В силу соотношения (2), имеем

$$p_n \sim \sum_{m \geq N} \frac{a}{m} \left(1 - \left(1 - \frac{x}{m^\beta}\right)^n\right) := \tilde{p}_n. \quad (3)$$

Рассмотрим близкую к \tilde{p}_n сумму ряда

$$s_n := \sum_{m \geq N} \frac{a}{m} \left(1 - e^{-\frac{nx}{m^\beta}}\right)$$

и покажем, что

$$\lim_{n \rightarrow \infty} s_n = \frac{a I(x/r^\beta)}{\beta}. \quad (4)$$

Действительно, заменяя сумму интегралом и делая замену переменной $u := \frac{nx}{t^\beta}$, получим

$$\sum_{m \geq N} \frac{a}{m} \left(1 - e^{-\frac{nx}{m^\beta}}\right) \sim \int_N^\infty \frac{adt}{t} \left(1 - e^{-\frac{nx}{t^\beta}}\right) = \frac{a}{\beta} \int_0^{\frac{x}{r^\beta}} \frac{1 - e^{-u}}{u} du = \frac{a I(x/r^\beta)}{\beta}.$$

Остаётся показать, что

$$\lim_{n \rightarrow \infty} (\tilde{p}_n - s_n) = 0. \quad (5)$$

Для этого напомним несколько элементарных оценок. При достаточно малых y имеем $\log(1 - y) \geq -y - y^2$, откуда при всех натуральных n выполнено

$$(1 - y)^n \geq \exp(-ny - ny^2)$$

и если $y > 0$, то

$$0 \leq \exp(-ny) - (1 - y)^n \leq \exp(-ny) (1 - \exp(-ny^2)) \leq ny^2.$$

Применяя эту оценку к $y = \frac{x}{m^\beta}$ для $m \geq N$, получим

$$\begin{aligned} 0 &\leq \tilde{p}_n - s_n = \sum_{m \geq N} \frac{a}{m} \left(e^{-\frac{nx}{m^\beta}} - \left(1 - \frac{x}{m^\beta}\right)^n\right) \\ &\leq \sum_{m \geq N} \frac{a}{m} n \frac{x^2}{m^{2\beta}} \sim \frac{ax^2}{2\beta} n N^{-2\beta} = \frac{ax^2}{2\beta r^{2\beta}} n^{-1} \rightarrow 0. \end{aligned}$$

Теперь из соотношений (3), (4) и (5) следует

$$\lim_{n \rightarrow \infty} p_n = \frac{a I(x/r^\beta)}{\beta},$$

что и требовалось доказать. \square

Доказательство теоремы 2. С одной стороны, очевидно, что $Q_j \geq e^{bV_j}$, $Q^n \geq e^{bV^n}$. Поэтому, с учётом теоремы 1,

$$\mathbb{P}(\log Q^n \leq brn^\gamma) \leq \mathbb{P}(bV^n \leq brn^\gamma) = \mathbb{P}(V^n \leq rn^\gamma) \rightarrow F_{\alpha, \beta}(r).$$

С другой стороны,

$$Q_j = \sum_{m \in A \cap X} e^{bm} \leq \sum_{m=1}^{V_j} e^{bm} \leq B e^{bV_j},$$

где $B := \frac{e^b}{e^b - 1}$.

Фиксируем $\varepsilon \in (0, 1)$. Если $V_j \leq (1 - \varepsilon)rn^\gamma$, то $Q_j \leq B e^{b(1 - \varepsilon)rn^\gamma}$. При достаточно больших n имеем $B e^{-\varepsilon brn^\gamma} < 1$, так что $Q_j \leq e^{brn^\gamma}$.

Соответственно из $V^n \leq (1 - \varepsilon)rn^\gamma$ следует $Q^n \leq e^{brn^\gamma}$, $\log Q^n \leq brn^\gamma$. По теореме 1 получаем

$$\mathbb{P}(\log Q^n \leq brn^\gamma) \geq \mathbb{P}(V^n \leq (1 - \varepsilon)rn^\gamma) \rightarrow F_{\alpha, \beta}((1 - \varepsilon)r).$$

Наконец, устремляя $\varepsilon \searrow 0$, приходим к

$$\liminf_{n \rightarrow \infty} \mathbb{P}(\log Q^n \leq brn^\gamma) \geq F_{\alpha, \beta}(r). \quad \square$$

СПИСОК ЛИТЕРАТУРЫ

1. F. Auerbach, *Das Gesetz der Bevölkerungskonzentration*. In: Petermanns Geogr. Mitteilungen **59** (1913), 73–76.
2. J.-B. Estoup, *Gammes Sténographiques*, various editions, Paris, 1908–1916.
3. B. Hoffmann, M. Lifshits, Yu. Lifshits, D. Nowotka, *Maximal intersection queries in randomized input models*. — Theory Comput. Syst. **46**, No. 1 (2010), 104–119.
4. A. Lelu, *Jean-Baptiste Estoup and the origins of Zipf's law: a stenographer with a scientific mind (1868–1950)*. — Boletn de Estadstica e Investigación Operativa **30**, No. 1 (2014), 66–77.
5. C. D. Manning, H. Schütze, *Foundations of Statistical Natural Language Processing*, MIT Press, Cambridge, 1999.
6. M. Petruszewycz, *L'histoire de la loi d'Estoup–Zipf: documents*. — Math. Sci. Hum. **44** (1973), 41–56.

7. G. K. Zipf, *Selected Studies of the Principle of Relative Frequency in Language*. Harvard Univ. Press, 1932.
8. G. K. Zipf, *The Psycho-Biology of Language*, Mifflin, Houghton, 1935.
9. М. П. Бакулина, *Использование закона Ципфа для сжатия текстов*. — Дискретн. анализ и исслед. опер. **14**, No. 2 (2007), 3–13.
10. Я. Галамбош, *Асимптотическая теория экстремальных порядковых статистик*, Наука, М., 1984.
11. М. А. Лифшиц, И. М. Лялинов, *Вероятностные свойства множеств Ципфа и их максимальных пересечений*, в печати, 2022.
12. В. П. Маслов, Т. В. Маслова, *О законе Ципфа и ранговых распределениях в лингвистике и семиотике*. — Матем. заметки **80**, No. 5 (2006), 718–732.
13. В. П. Маслов, *Фазовые переходы нулевого рода и квантование закона Ципфа*. — Теор. матем. физ. **150**, No. 1 (2007), 118–142.
14. В. П. Маслов, *Об одной общей теореме теории множеств, приводящей к распределению Гиббса, Бозе–Эйнштейна, Парето и закону Ципфа–Мандельброта для фондового рынка*. — Матем. заметки **78**, No. 6 (2005), 870–877.
15. Ю. И. Манин, *Закон Ципфа и вероятностные распределения Левина*. — Функц. анализ и его прил. **48**, No. 2 (2004), 51–66.
16. Ю. А. Шрейдер, *О возможности теоретического вывода статистических закономерностей текста (к обоснованию закона Ципфа)*. — Пробл. передачи информ. **3**, No. 1 (1967), 57–63.
17. Б. А. Трубников, И. А. Румынский, *Простейший вывод закона Ципфа–Крылова для слов и возможность его “эволюционной” интерпретации*. — Докл. АН СССР **321**, No. 2 (1991), 270–275.

Lifshits M. A., Lyalinov I. M. Two limit theorems on the intersections of random Zipf sets.

In this work we study the asymptotic behaviour of the rarest element in the intersections of a random Zipf set with a large number of independent random sets of the same type but, eventually, with different parameters. The same problem is solved for the maximum of the integral measure of intersection associated with exponentially growing weights.

С.-Петербургский
государственный университет,
Университетская наб. 7/9,
191023, Санкт-Петербург,
Россия

E-mail: mikhail@lifshits.org

E-mail: lyalinov239@yandex.ru

Поступило 29 августа 2022 г.