

*На правах рукописи*

**Тимофеева Нина Евгеньевна**

**ПОСТРОЕНИЕ ОЦЕНОК ЭНТРОПИИ  
СТАЦИОНАРНЫХ СЛУЧАЙНЫХ ПРОЦЕССОВ**

Специальность 05.13.18 - математическое моделирование,  
численные методы и комплексы программ

**АВТОРЕФЕРАТ**

диссертации на соискание ученой степени  
кандидата физико-математических наук

Ярославль – 2010

Работа выполнена на кафедре теории функций и функционального анализа Ярославского государственного университета им. П.Г. Демидова

**Научный руководитель** – доктор физико-математических наук, профессор Дольников Владимир Леонидович

**Официальные оппоненты:** доктор физико-математических наук, профессор Соколов Валерий Анатольевич  
доктор физико-математических наук, доцент Райгородский Андрей Михайлович

**Ведущая организация** – Учреждение Российской академии наук  
Институт программных систем имени А.К.Айламазяна РАН

Защита состоится "\_\_\_" \_\_\_\_\_ 2010 г. в \_\_\_ часов на заседании диссертационного совета Д 212.002.05 при Ярославском государственном университете имени П.Г. Демидова по адресу: 150000, г. Ярославль, ул. Советская, д. 14.

С диссертацией можно ознакомиться в библиотеке Ярославского государственного университета имени П.Г. Демидова по адресу: 150000, г. Ярославль, ул. Полушкина роща, д. 1.

Автореферат разослан "\_\_\_" \_\_\_\_\_ 2010 г.

Ученый секретарь диссертационного совета

Глызин С.Д.

## ОБЩАЯ ХАРАКТЕРИСТИКА РАБОТЫ

### **Актуальность темы.**

Математическое понятие энтропии впервые ввел К.Шеннон в 1948 г. Он использовал его для измерения количества информации, показал, что энтропия определяет границу степени сжатия текста.

В 50-е гг. А.Н. Колмогоров перенес понятие энтропии на динамические системы и доказал, что энтропия есть инвариант динамической системы: изоморфные динамические системы имеют одинаковую энтропию ([11]). Он применил понятие энтропии для решения задачи об изоморфности сдвигов Бернулли.

В 1971 г. Д.Орнштейн доказал утверждение, обратное теореме Колмогорова: если энтропии сдвигов Бернулли одинаковы, то эти сдвиги изоморфны[12].

Эти исследования энтропии имеют важное теоретическое значение для математики, в частности, для информатики.

В прикладных задачах возникает проблема нахождения значения энтропии, поскольку она определяет границу сжатия информации. По экспериментальным данным можно только ценить энтропию. Существуют два основных подхода к построению оценок энтропии: построение эмпирической функции распределения и непараметрические оценки.

В первом подходе на прямую вычисляется математическое ожидание логарифма эмпирической функции распределения. Этот способ удобен, если число параметров, от которых зависит распределение, конечно. Впервые такой подход описан Г.П.Башариным[10].

Непараметрические оценки энтропии могут быть поделены на два больших класса. Первый класс использует алгоритмы сжатия данных Лемпеля-Зива[1]. Основываясь на результате Лемпеля-Зива, П.Грасбергер предложил свою первую оценку величины обратной к энтропии[2]. Используя результаты Д.С.Орнштейна и Б.Вейса, П.Шилдс доказал, что эта оценка не является сходящейся для общих эргодических процессов, но при наложении определенных условий будет состоятельной [9]. Основываясь на этом результате, И.Контояннис и Ю.М.Сухов показали состоятельность оценки для

более широкого класса стационарных эргодических процессов [3].

Второй класс оценок основан на "методе расстояния до ближайших соседей". Р.Л.Добрушин первым предложил оценку энтропии, использующую этот метод.[5]. Основная идея заключается в следующем: если ранее рассматривалась одна бесконечная последовательность, то при новом подходе рассматривается бесконечное число конечных последовательностей и вычисляется расстояние между ними. В.А.Ватутин и В.Г.Михайлов нашли смещение оценки Р.Л.Добрушина, оценили значение ее дисперсии и доказали состоятельность [6]. Позднее исследователи перенесли идею Р.Л.Добрушина на пространство случайных последовательностей. Опираясь на полученные результаты, Р.Бадии и А.Полити. предложили оценивать размерность Хаусдорфа в общем метрическом пространстве [7]. П.Биллингслей показал, что при выборе подходящей метрики энтропия совпадает с размерностью Хаусдорфа [8]. Основываясь на результатах Р.Бадии и А.Полити, П.Грассбергер ввел свою вторую оценку энтропии на пространстве случайных последовательностей [2]. П.Шилдс построил пример, который показывает, что оценка П.Грассбергера несостоятельна для общего эргодического процесса. Также он показал, что предлагаемая П.Грассбергером оценка состоятельна для неприводимых непериодических марковских цепей [9]. В работе В.В.Майорова и Е.А.Тимофеева было введено обобщение оценки П.Грассбергера[13].

Таким образом, теория энтропии является одной из активно развивающихся областей математики и находит применение во многих областях современной науки.

**Цель работы.** Построение новых статистических оценок энтропии. Исследование свойств этих оценок, а именно: смещенности, состоятельности, порядка убывания дисперсии. Построение эффективного алгоритма для вычисления энтропии динамических систем по их траекториям.

**Методы исследования.** В работе используются методы теории вероятности, теории динамических систем, теории информации, линейной алгебры.

**Научная новизна.** Основные результаты диссертации являются

новыми и состоят в следующем.

Построены три новые статистические оценки энтропии, основанные на методе "расстояния до ближайших соседей". Принципиальные преимущества новых оценок это: степенной порядок точности дисперсии, который уменьшен почти до границы Рао-Крамера; при оценивании энтропии динамических систем значения оценок не зависят от разбиения пространства (это подтверждается экспериментальными исследованиями); в ряде случаев показан степенной порядок убывания дисперсии.

Исследованы свойства этих оценок: смещенность, состоятельность. Доказано, что первая из предложенных оценок сходится  $\mu$ -почти всюду. Описан прием уменьшения смещения оценки. Проведены вычислительные эксперименты, подтверждающие точность предложенных оценок.

В работе также приведен эффективный алгоритм вычисления оценки по экспериментальным данным.

**Теоретическая и практическая ценность.** Работа носит теоретический характер. Результаты диссертации могут быть полезны при изучении задач построения оценок энтропии, а также при изучении динамических систем.

**Апробация диссертации.** Основные диссертации докладывались на XI симпозиуме по проблемам избыточности в информационных и управляющих системах, г. Санкт-Петербург, 2007 г., на 45 ежегодной конференции университета штата Иллинойс США, 2007 г., на IEEE Information Theory Workshop, Lake Tahoe, California, USA, 2007 г., на Международной конференции "Компьютерные технологии в электротехнике и радиоэлектронике", Новосибирск, 2008г., на семинаре кафедры математического моделирования Ярославского государственного университета им. П.Г. Демидова.

**Публикации.** По теме диссертации опубликовано 10 научных работ, из них 5 в журналах из перечня ВАК. Список публикаций приведен в конце автореферата. Работы, написанные совместно с другими исследователями, выполнены в нераздельном соавторстве.

**Структура диссертации.** Диссертация состоит из введения и 4-х глав. Список литературы содержит 39 наименований. Общий объем

диссертации составляет 63 страницы.

## СОДЕРЖАНИЕ ДИССЕРТАЦИИ

Рассматривается пространство последовательностей  $\Omega = \mathcal{A}^{\mathbb{N}}$ . Пусть  $\mu$  — эргодическая, инвариантная относительно сдвига мера на пространстве  $\Omega$ . Пусть  $\xi_1, \dots, \xi_n$  — независимые, одинаково распределенные по мере  $\mu$  случайные точки из  $\Omega$ .

В прикладных задачах задано  $n$  слов длины  $m$ . В этом случае, считаем, что эти слова совпадают с первыми  $m$  символами точек  $\xi_1, \dots, \xi_n$  и будем обозначать эти слова через  $\xi_1^{(m)}, \dots, \xi_n^{(m)}$ .

Требуется построить оценку величины обратной к энтропии. Эта замена энтропии на обратную величину сделана для удобства формулировки и доказательства теорем.

### Глава 1.

В главе 1 на пространстве  $\Omega$  рассматривается метрика

$$\rho_0(\mathbf{x}, \mathbf{y}) = \frac{1}{\min\{k : x_k \neq y_k\}}. \quad (1)$$

Предлагается следующая оценка:

$$\eta_n^{(k,m)} = \frac{1}{n \log n} \sum_{j=1}^n \left( \min_{i:i \neq j}^{(k)} \rho^{(m)}(\xi_i, \xi_j) \right)^{-1}, \quad (2)$$

где

$$\rho_0^{(m)}(\mathbf{x}, \mathbf{y}) = \max \left\{ \frac{1}{m}, \rho_0(\mathbf{x}, \mathbf{y}) \right\}. \quad (3)$$

**Замечание.** Подчеркнем, что предлагаемая оценка использует только первые  $m$  символов последовательностей  $\xi_1, \dots, \xi_n$ .

Асимптотическая несмещенность оценки доказана в следующей теореме.

**Теорема 1.** Пусть  $\mu$  — инвариантная относительно сдвига борелевская вероятностная эргодическая мера, удовлетворяющая условию:

$$\exists a > 0, b : \mu(C_s(\mathbf{x})) \leq e^{-as+b}, \quad \forall s > 0, \mathbf{x} \in \Omega \quad (4)$$

Тогда для всех  $m \geq \frac{2}{a} \ln n$  существует предел:

$$\lim_{n \rightarrow \infty} \mathbb{E} \eta_n^{(k,m)} = \frac{1}{h}. \quad (5)$$

где  $\eta_n^{(k,m)}$  определена в (2).

Отметим, что условия теоремы выполняются для гиббсовских мер.

По сравнению с ранее известным результатом П.Шилдса сходимость оценки (2) показана для более широкого класса мер.

Для любого  $s \in \mathbb{R}$  и произвольной вещественной функции  $\phi : \mathbb{R} \rightarrow \mathbb{R}$ , определим ее конечную разность  $\Delta\phi(s)$  как

$$\Delta\phi(s) = \phi(s) - \phi(s-1) \quad (6)$$

и  $k$ -ю разность как

$$\Delta^{(k)}\phi(s) = \Delta^{(k-1)}\phi(s) - \Delta^{(k-1)}\phi(s-1). \quad (7)$$

Обозначим:

$$\begin{aligned} \phi_m(n) &= \sum_{s=0}^{m-1} \left( 1 - \int_{\Omega} [1 - \mu(C_s(x))]^n d\mu(x) \right), \\ \phi(n) &= \sum_{s=0}^{\infty} \left( 1 - \int_{\Omega} [1 - \mu(C_s(x))]^n d\mu(x) \right). \end{aligned} \quad (8)$$

**Лемма 1.**

$$\mathbb{E} \eta_{n+1}^{(k,m)} = \frac{1}{\log n} \sum_{j=0}^{k-1} (-1)^j \binom{n}{j} \Delta^{(j)}\phi_m(n). \quad (9)$$

**Лемма 2.** Пусть мера  $\mu$  удовлетворяет условию (4). Тогда, для  $m \geq \frac{2}{a} \ln n$  и заданного  $k$ , существует такая константа  $C$ , что следующее неравенство выполнено:

$$\mathbb{E} \eta_{n+1}^{(k,\infty)} - \mathbb{E} \eta_{n+1}^{(k,m)} \leq Cn^{-1}.$$

Неравенство для параметра  $m$ , полученное в этой лемме, ранее было доказано П.Шилдсом для мер с условием свободной энергии. Условие леммы 2 является более широким.

В следующей теореме найдено ограничение сверху для дисперсии оценки.

**Теорема 2.** Пусть  $\mu$  удовлетворяет условиям теоремы 1. Тогда выполнено следующее неравенство:

$$D\eta_n^{(k,m)} \leq \frac{(2k-1)^2(m-1)^2}{4n(\log n)^2} \quad (10)$$

где  $\eta_n^{(k,m)}$  определена в (2).

**Замечание.** Для  $m = \mathcal{O}(\log n)$  и фиксированного параметра  $k$  оценка (10) равняется  $\mathcal{O}(n^{-1})$ .

Полученное неравенство для дисперсии в теореме 2 при  $m = \mathcal{O}(\log n)$ , улучшает ранее известный результат —  $\mathcal{O}(n^{-c})$ . Более того, полученный порядок убывания оценки по параметру  $n$  совпадает с порядком классической нижней границы Рао-Крамера для дисперсии.

Таким образом, предложенная в главе 1 оценка является наилучшей в настоящий момент по порядку убывания дисперсии по  $n$  и этот порядок, скорее всего, не улучшаем.

Из теорем 1 и 2 легко выводится состоятельность оценки.

Сходимость почти всюду была доказана П.Шилдсом [9] для марковских мер, а в данной работе получена сходимость для более широкого класса мер.

**Теорема 3.** Пусть мера  $\mu$  удовлетворяет условиям теоремы 1. Тогда оценка  $\eta_n^{(k,m)}$ , определенная в (2), сходится к  $\frac{1}{h}$   $\mu$ -почти всюду.

В работе проведен эксперимент для нескольких динамических систем, энтропия которых известна. Результаты эксперимента подтверждают эффективность оценки (2) для нахождения энтропии.



**Глава 2.** На пространстве  $\Omega$  рассматривается метрика

$$\rho_1(\mathbf{x}, \mathbf{y}) = \frac{\theta - 1}{|\mathcal{A}| - 1} \left| \sum_{k=1}^{\infty} \theta^{-k} x_k - \sum_{k=1}^{\infty} \theta^{-k} y_k \right|, \quad (11)$$

где параметр  $\theta \geq |\mathcal{A}|$ ,  $|\mathcal{A}|$  — число символов алфавита  $\mathcal{A}$ .

Построена следующая оценка энтропии:

$$\eta_n^{(k,m)} = \frac{k \left( r_n^{(k,m)} - r_n^{(k+1,m)} \right)}{\ln \theta}, \quad (12)$$

где

$$r_n^{(k,m)} = -\frac{1}{(n+1)} \sum_{j=1}^n \ln \left( \min_{i:i \neq j}^{(k)} \rho_1^{(m)}(\xi_i, \xi_j) \right), \quad (13)$$

метрика  $\rho_1$  определена по формуле (11) и через  $\rho_1^{(m)}(\mathbf{x}, \mathbf{y})$  обозначается усечение метрики:

$$\rho_1^{(m)}(\mathbf{x}, \mathbf{y}) = \max \left\{ \frac{\theta - 1}{|\mathcal{A}| - 1} \theta^{-m}, \rho_1(\mathbf{x}, \mathbf{y}) \right\}. \quad (14)$$

Так как

$$\lim_{\theta \rightarrow \infty} \frac{\ln \rho_1(\mathbf{x}, \mathbf{y})}{\ln \theta} = \rho_0(\mathbf{x}, \mathbf{y})^{-1},$$

то метрику  $\rho_0$  можно рассматривать как частный случай метрики  $\rho_1$  при  $\theta = \infty$ . Поэтому при  $\theta \rightarrow \infty$  оценка (12) совпадает с оценкой, рассмотренной в главе 1.

Подчеркнем, что предлагаемая оценка (12) использует только  $m$  первых координат заданных случайных точек  $\Omega$ . Эта естественная модификация оценки из работы [14] позволила улучшить оценку дисперсии.

Через  $\lambda$  обозначим отображение  $\lambda : \Omega \rightarrow [0, 1]$

$$\lambda(\omega) = \frac{\theta - 1}{|\mathcal{A}| - 1} \sum_{k=1}^{\infty} \theta^{-k} \omega_k. \quad (15)$$

При  $\theta > |\mathcal{A}|$  это отображение будет вложением и изоморфизмом компакта  $\Omega$  и компакта  $\lambda(\Omega) \subset [0, 1]$ . При  $\theta = |\mathcal{A}|$  это отображение будет метрическим изоморфизмом.

Отображение  $\lambda$  переносит меру  $\mu$  на  $\lambda(\Omega)$ . Распространим полученную меру на весь отрезок  $[0, 1]$  и обозначим через  $\widehat{\mu}(x)$  соответствующую меру отрезка  $[0, x]$  (функцию распределения). Будем считать, что  $\widehat{\mu}(x) = 0$  при  $x < 0$ , и  $\widehat{\mu}(x) = 1$  при  $x > 1$ .

Пусть

$$\beta(x, r) = \widehat{\mu}(x + r) - \widehat{\mu}(x - r), \quad 0 \leq r \leq 1, \quad -1 \leq x \leq 2. \quad (16)$$

Ясно, что

$$\beta(x, r) = \mu(\{\omega \in \Omega : |\lambda(\omega) - x| < r\}), \quad 0 \leq r \leq 1, \quad -1 \leq x \leq 2. \quad (17)$$

Через  $r = \nu(x, t)$  обозначим обратную (по  $r$ ) функцию к функции  $t = \beta(x, r)$ .

Положим

$$\chi(t) = - \int_0^1 \ln \nu(x, t) d\widehat{\mu}(x). \quad (18)$$

Через  $\preceq$  будем обозначать лексикографический порядок на пространстве  $\Omega$ .

Нетрудно видеть, что метрика  $\rho_1$  согласована с порядком, т.е.

$$\rho_1(\mathbf{x}, \mathbf{y}) \leq \rho_1(\mathbf{x}, \mathbf{z}) \quad \forall \mathbf{x} \preceq \mathbf{y} \preceq \mathbf{z}. \quad (19)$$

Обозначим  $E r_n^{(k)} = E r_n^{(k, \infty)}$ ,  $E \eta_n^{(k)} = E \eta_n^{(k, \infty)}$ .

**Лемма 3.**

$$E r_n^{(k)} = \int_0^1 \int_0^1 \left[ 1 - \sum_{j=0}^{k-1} \binom{n}{j} \beta(x, u)^j (1 - \beta(x, u))^{n-j} \right] d\widehat{\mu}(x) \frac{du}{u}. \quad (20)$$

**Лемма 4.**

$$\mathbb{E}r_n^{(k,m)} = \int_{\varepsilon_m}^1 \int_0^1 \left[ 1 - \sum_{j=0}^{k-1} \binom{n}{j} \beta(x, u)^j (1 - \beta(x, u))^{n-j} \right] d\widehat{\mu}(x) \frac{du}{u}, \quad (21)$$

где

$$\varepsilon_m = \frac{\theta - 1}{|\mathcal{A}| - 1} \theta^{-m}.$$

**Лемма 5.** Пусть мера  $\mu$  удовлетворяет условию

$$\lim_{t \rightarrow 0} t \ln \nu(t, x) = 0, \quad \forall x, \quad (22)$$

тогда

$$\mathbb{E}r_n^{(k)} = k \binom{n}{k} \int_0^1 \chi(t) t^{k-1} (1-t)^{n-k} dt. \quad (23)$$

$$\mathbb{E}\eta_n^{(k)} = \frac{k}{\ln \theta} \binom{n}{k} \int_0^1 \int_0^1 \beta(x, u)^k (1 - \beta(x, u))^{n-k} d\widehat{\mu}(x) \frac{du}{u}. \quad (24)$$

**Теорема 4.** Пусть  $\mu$  – инвариантная относительно сдвига борелевская вероятностная мера и  $\xi_0, \dots, \xi_n$  – независимые случайные точки в  $\Omega$  распределенные по мере  $\mu$ , тогда справедливо неравенство

$$\mathbb{D}r_n^{(k,m)} \leq \frac{(2k-1)^2 m^2}{4(n+1)} \ln^2 \theta, \quad (25)$$

где  $r_n^{(k,m)}$  определено в (13).

**Следствие 1.** Пусть  $\mu$  – инвариантная относительно сдвига борелевская вероятностная мера и  $\xi_0, \dots, \xi_n$  – независимые случайные точки в  $\Omega$  распределенные по мере  $\mu$ , тогда справедливо неравенство

$$\mathbb{D}\eta_n^{(k,m)} \leq \frac{4k^4 m^2}{n+1}. \quad (26)$$

Правило выбора параметра  $m$  установлено в следующих утверждениях.

**Лемма 6.** Пусть  $\mu$  – инвариантная относительно сдвига борелевская вероятностная мера, удовлетворяющая условию

$$\exists a > 0, C > 0 : \beta(x, t) \leq Ct^a, \quad \forall t > 0, x \in [0, 1]. \quad (27)$$

Тогда существует константа  $C_1$  такая, что:

$$\mathbb{E} r_n^{(k, \infty)} - \mathbb{E} r_n^{(k, m)} \leq C_1 n^k \theta^{-mk}.$$

**Следствие 2.** Пусть  $\mu$  – инвариантная относительно сдвига борелевская вероятностная мера, удовлетворяющая условию (27), тогда для  $m \geq \frac{2}{a \ln \theta} \ln n$  и фиксированного  $k$ , существует константа  $C_2$  такая, что:

$$\mathbb{E} \eta_n^{(k, \infty)} - \mathbb{E} \eta_n^{(k, m)} \leq C_2 n^{-1}.$$

Доказана сходимость оценки (12) почти всюду.

**Теорема 5.** Пусть  $\mu$  – инвариантная относительно сдвига борелевская вероятностная мера, удовлетворяющая условию (27), тогда,

$$\lim_{n \rightarrow \infty} \left[ \eta_n^{(k, m)} - \mathbb{E} \eta_n^{(k, \infty)} \right] = 0 \text{ почти всюду}$$

при  $\frac{2}{a \ln \theta} \ln n < m \leq C \ln n$ , где  $C$  – некоторая константа.

Для симметричной меры Бернулли смещение оценки (12) является степенным по  $n$ .

$$\mathbb{E} \eta_n^{(k, \infty)} = \frac{1}{\ln |\mathcal{A}|} \left( 1 - \frac{2 \ln 2 - 1}{n + 1} \right), \quad (28)$$

т.е. при  $n \rightarrow \infty$  математическое ожидание оценки стремится к  $1/h$  ( $h = |\ln \mathcal{A}|$  для симметрической меры Бернулли).

Итак, скорость сходимости оценки (12) для симметричной меры Бернулли при  $\theta = |\mathcal{A}|$  равна  $\mathcal{O}(n^{-1})$  при  $m = \mathcal{O}(\ln n)$ .

Таким образом, предложенная в главе 2 оценка энтропии для этой меры лучше, чем оценка из первой главы, у которой смещение равно  $\mathcal{O}(1)$ . Ее недостаток, по сравнению с оценкой (2) первой главы, — не доказан факт сходимости к  $1/h$  почти всюду.

### Глава 3.

В этой главе описан общий подход к построению эвристических оценок энтропии и исследованы их свойства.

Опишем основную идею предлагаемого подхода. Исследованная в главе 1 оценка содержит целочисленный параметр  $k$ . Новые оценки энтропии строятся как линейная комбинация этих оценок при различных значениях параметра  $k$ . Коэффициенты линейной комбинации подбираются так, чтобы формальный ряд по  $n$  смещения оценки имел заданное число нулевых случайных слагаемых.

Этот подход справедлив не только для метрики  $\rho_0$ , но и для произвольной метрики  $\rho$ , поэтому рассматривается произвольная метрика  $\rho$  на пространстве последовательностей.

Исследуемая оценка задается формулой

$$\eta_n^{(k,m)} = \sum_{i=1}^k a_i r_n^{(i,m)}, \quad (29)$$

где

$$a_i \triangleq (-1)^{i-1} (k-1) \binom{k-1}{i-1}, \quad (30)$$

а  $r_n^{(k,m)}$  определена в формуле (13) и для простоты будем считать, что  $m = \infty$ . Для краткости будем обозначать  $\eta_n^{(k)} = \eta_n^{(k,m)}$ ,  $r_n^{(k)} = r_n^{(k,m)}$ .

Выбор коэффициентов  $a_i$  объясняется в следующем утверждении.

**Утверждение 1.** Пусть величина  $E\eta_n^{(k)}$  представлена формальным рядом по  $n^{-1}$  и коэффициенты  $a_i$  в формуле (29) определены согласно формуле (30). Тогда формальный ряд для  $E\eta_n^{(k)}$  имеет вид:

$$E\eta_n^{(k)} = \frac{1}{h} + \psi_{k-1} n^{-k+1} + \psi_k n^{-k} + \dots$$

Положим

$$\phi(n) = \text{Er}_n^{(1)}. \quad (31)$$

Сходимость оценки (29) показана в следующей теореме.

**Теорема 6.** *Предположим, что существуют константы  $h > 0$  и  $a > 0$  такие, что для некоторого  $k = K_0 > 0$  верно*

$$\Delta^{(k)}\phi(n) = (-1)^{k-1} \frac{(k-1)!(n-k)!}{hn!} (1 + \mathcal{O}(n^{-a})). \quad (32)$$

Тогда константа  $h$  из (32) равна энтропии меры  $\mu$  для всех  $k < K_0$  и верно следующее равенство:

$$\lim_{n \rightarrow \infty} \text{Er}_n^{(k)} = \frac{1}{h}.$$

Доказательство теоремы основано на двух, интересных на наш взгляд, леммах.

**Лемма 7.** *Если условие (32) выполнено для некоторого  $k = k_0 > 0$ , то оно выполнено также для  $0 < k \leq k_0$  и*

$$\phi(n) = \frac{H_n}{h} + C_0 + \mathcal{O}(n^{-a}), \quad (33)$$

где

$$C_0 = \sum_{s=1}^{\infty} \left( \Delta\phi(s) - \frac{1}{hs} \right).$$

**Лемма 8.** *Предположим, что существуют константы  $h > 0$ ,  $a > 0$  и  $C > 0$  такие, что*

$$\Delta^{(k)}\phi(n) = \frac{(-1)^{k-1}}{hk} \binom{n}{k}^{-1} (1 + \psi_k(n)), \quad (34)$$

выполнено для некоторого  $k = k_0 > 0$ , где

$$|\psi_k(n)| \leq Cn^{-a}. \quad (35)$$

Тогда условия (34) и (35) верны для всех  $0 < k \leq k_0$  с некоторой константой  $C$  и справедливы следующие равенства:

$$\phi(n) = \frac{H_n}{h} + C_0 + \psi_0(n), \quad (36)$$

где  $\psi_0(n)$  определена в (35) и  $C_0$  определена в лемме 7.

Вычислительный эксперимент подтвердил эффективность предложенного подхода.

#### Глава 4.

В этой главе описан экономный (по трудоемкости) алгоритм, который позволяет наблюдать смещение оценки. При вычислении исследуемых оценок энтропии и построения оценок размерности фракталов в работе [13] находят следующую величину:

$$r_n^{(k)} = \frac{1}{n+1} \sum_{j=0}^n \phi \left( \min_{i:i \neq j}^{(k)} \rho(\xi_i, \xi_j) \right), \quad (37)$$

где  $\phi$  — некоторая монотонная функция.

Предлагаемый в диссертации алгоритм позволяет находить величины  $r_s^{(k)}$  для всех значений  $s \leq n$  с такой же трудоемкостью (по порядку роста) как и величину  $r_s^{(k)}$ .

Трудоемкость алгоритма равна  $\mathcal{O}(mn)$ , где  $m$  — длина слов. Нахождение оценок для различных значений  $n$  позволяет визуально наблюдать смещение, которое может зависеть от  $n$ .

#### Постановка задачи.

Пусть заданы  $n+1$  последовательностей  $\xi_0 = (x_{01}, \dots, x_{0m}), \dots, \xi_n = (x_{n1}, \dots, x_{nm})$ , где  $x_{ij} \in \mathcal{A} = \{1, 2, \dots, a\}$ .

Будем предполагать, что применяемая метрика  $\rho$  на  $\Omega$  она согласована с лексикографическим порядком, т.е.

$$\rho(\mathbf{x}, \mathbf{y}) \leq \rho(\mathbf{x}, \mathbf{z}) \quad \forall \mathbf{x} \preceq \mathbf{y} \preceq \mathbf{z}. \quad (38)$$

Алгоритм находит величины  $r_s^{(l)}$ , определенные по формуле (37), для всех  $k < s \leq n$ ,  $1 \leq l \leq k$ , где  $k$  — заданное число.

#### Описание алгоритма.

Работа алгоритма состоит из трех шагов.

1. Выполним сортировку слов  $\xi_0, \dots, \xi_n$  в возрастающем порядке.  
Пусть

$$\eta_0 \preceq \dots \preceq \eta_n$$

– отсортированные слова, а  $\sigma$  – перестановка такая, что

$$\eta_{\sigma_i} = \xi_i, \quad i = 0, 1, \dots, n.$$

Организуем слова  $\eta_0, \dots, \eta_n$  в список с двумя связями.

2. Для  $j = 0, 1, \dots, n$  и  $l = 1, 2, \dots, k$  находим

$$g_j^l = \min_{i \neq j: j-l \leq i \leq j+l} \binom{k}{i} \rho(\eta_i, \eta_j). \quad (39)$$

Затем находим

$$r_n^{(l)} = \frac{1}{n+1} \sum_{j=0}^n \phi(g_j^l), \quad (40)$$

где  $l = 1, 2, \dots, k$ .

3. Для  $s = n, n-1, \dots, k$  делаем следующее:

- (а) удаляем слова  $\eta_{\sigma_s} = \xi_s$  из списка,
- (б) для  $l = 1, 2, \dots, k$  и  $j = \sigma_s - l, \dots, \sigma_s + l$  находим величины  $g_j^l$  по формуле (39); нумерация по  $j$  идет по списку с двумя связями;
- (с) для  $l = 1, 2, \dots, k$  пересчитываем величины  $r_{s-1}^{(l)}$ , заменяя в формуле (40) величины  $g_j^l$  для  $j = \sigma_s - l, \dots, \sigma_s + l$ .

Отметим, что величины  $g_j^l$  находятся последовательно для  $l = 1, 2, \dots, k$ , поэтому для нахождения каждой величины  $g_j^l$  по формуле (39) достаточно находить минимум только из двух расстояний.

Легко видеть, что трудоемкость алгоритма равна  $\mathcal{O}(k^2nm)$ . Таким образом, при небольших  $k$  (не зависящих от  $n$ ) трудоемкость алгоритма равна  $\mathcal{O}(nm)$ .



Свойство согласованности (38) выполнено для метрик  $\rho_0$  и  $\rho_1$  (см. формулы (1) и (11) соответственно), однако подчеркнем, что оно не выполнено для метрики

$$\rho_2(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^m \theta^{-k} |x_k - y_k|,$$

где  $\theta > 1$ .

## Список литературы

1. A. Wyner and J. Ziv. Some asymptotic properties of the entropy of a stationary ergodic data source with applications to data compression. / A. Wyner and J. Ziv. - IEEE Trans. Inform. Theory, 1989. - No 35. - pp.1250–1258.
2. Grassberger, P. Estimating the information content of symbol sequences and efficient codes. /P. Grassberger - IEE Trans. Inform. Theory, 1989. - No 35. - pp. 669-675.
3. I. Kontoyiannis and Yu. M. Suhov. Prefixes and the entropy rate for long-range sources / I. Kontoyiannis and Yu. M. Suhov. - Probability Statistics and Optimization; ed. F.P. Kelly. - Wiley, New York, 1994. - pp. 89-98.
4. I. Kontoyiannis, P. H. Algoet, Yu. M. Suhov and A. J. Wyner. Nonparametric entropy estimation for stationary processes and random fields, with applications to English text. /I. Kontoyiannis, P. H. Algoet, Yu. M. Suhov and A. J. Wyner. - IEEE Trans. Inform. Theory, 1998. - N 44. - pp. 1319-1327.
5. Dobrushin, R.L. The simplified method of experimental estimation of the entropy of stationary sequence. / R.L. Dobrushin. - Theor. Prob. Appl., 1958, - No 3. - pp. 462-464.

6. V. A. Vatutin and V. G. Mikhaïlov. Statistical estimation of the entropy of discrete random variables with a large number of outcomes. /V. A. Vatutin and V. G. Mikhaïlov. - Uspekhi Mat. Nauk, 1995. - No 55. - pp. 121-134.
7. R. Badii and A. Politi. Hausdorff dimension and uniformity factor of strange attractors. /R. Badii and A. Politi. - Phys. Rev. Lett., 1984. - No 55. - pp. 1661-1664.
8. Billingsley, P. Ergodic Theory and Information. / P.Billingsley. - Wiley, New York, 1965. - pp. 120.
9. Shields, P.C.. Entropy and prefixes. /P.C. Shields. - Annals of Probability, 1992. - No 20. - pp. 403-409.
10. Башарин, Г.П. О статистическом оценивании энтропии последовательности независимых случайных величин. /Г.П. Башарин. Теория вероятности и ее применение, 1959. - Т. IV, No 3. - С.361-364.
11. Колмогоров, А.Н. Теория информации и теория алгоритмов./ А.Н. Колмогоров. - М.: Наука, 1987. - 300 с.
12. Орнстейн, Д. Эргодическая теория, случайность и динамические системы. / Д. Орнстейн. - М.: Мир, 1971. - 166 с.
13. В.В. Майоров, Е.А.Тимофеев. Статистическая оценка обобщенных размерностей. /В.В. Майоров, Е.А.Тимофеев. // Мат. заметки. - 2002. - Т. 71. № 5. - С.679 - 712.
14. Тимофеев, Е.А. Статистически оцениваемые инварианты мер. / Е.А. Тимофеев // Алгебра и анализ. - 2005. - Т.17. No 3. - С.204-236.

## РАБОТЫ АВТОРА ПО ТЕМЕ ДИССЕРТАЦИИ

### **Статьи в журналах из перечня ВАК:**

1. Тимофеева, Н.Е. Экономный алгоритм нахождения средних минимальных расстояний / Н.Е. Тимофеева // Моделирование и анализ информационных систем / Яросл. гос. ун-т. - Ярославль: ЯрГУ, 2007. - Т.14, № 3. - С.50-52.

2. A. Kaltchenko, N. Timofeeva. Entropy Estimators with Almost Sure Convergence and an  $O(n^{-1})$  Variance / A. Kaltchenko and N. Timofeeva // Advances in Mathematics of Communications. - 2008. - Vol. 2, No 1. - pp. 1–13

3. A. Kaltchenko, N. Timofeeva, and E. Timofeev. Bias Reduction of the Nearest Neighbor Entropy Estimator / A. Kaltchenko, N. Timofeeva, and E. Timofeev // International Journal of Bifurcation and Chaos. - 2008. - pp.3781–3787.

4. Тимофеева, Н.Е. Линейная метрика для оценивания энтропии / Н.Е. Тимофеева // Моделирование и анализ информационных систем / Яросл. гос. ун-т. - Ярославль: ЯрГУ, 2009. - Т.16, № 1. - С.39-48.

5. A. Kaltchenko and N. Timofeeva. Rate of convergence of the Nearest Neighbor Entropy Estimator / A. Kaltchenko and N. Timofeeva // International Journal of Electronics and Communications. - 2010. - Vol.64. - pp.75–79.

### **Работы, опубликованные в других журналах:**

6. A. Kaltchenko, E-H. Yang, and N. Timofeeva. Entropy Estimators with Almost Sure Convergence and an  $O(n^{-1})$  Variance / A. Kaltchenko, E-H. Yang, and N. Timofeeva // Proceedings of the 2007 IEEE Information Theory Workshop / Lake Tahoe, California, USA, 2007. - pp.644 - 649.

7. A. Kaltchenko, I. Kotsireas, E-H. Yang, and N. Timofeeva. Entropy Rate Estimators with a Near-Optimal Upper Bound on Variance / A. Kaltchenko, I. Kotsireas, E-H. Yang, and N. Timofeeva // Proceedings of the XI International Symposium on Problems of Redundancy in Information and Control Systems / Saint Petersburg, Russia, 2007. - pp. 18 – 21.

8. A. Kaltchenko, E-H. Yang, and N. Timofeeva. Exact Analysis of the Bias of the Nearest Neighbor Entropy Estimator for I.I.D. Information Sources / A. Kaltchenko, E-H. Yang, and N. Timofeeva // Proceedings of the Forty-Fifth Annual Allerton Conference / University of Illinois at Urbana-Champaign, IL, USA, 2007. - pp. 165– 168.

9. A. Kaltchenko, E-H. Yang, and N. Timofeeva. Bias Reduction of the Nearest Neighbor Entropy Estimator / A. Kaltchenko, E-H. Yang, and N. Timofeeva // Proceedings of the 2008 IEEE International Conference on Computational Technologies in Electrical and Electronics Engineering / Novosibirsk, Russia, 2008. - pp. 261 – 265.

10. A. Kaltchenko, N. Timofeeva. Bias Reduction via Linear Combination of Nearest Neighbor Entropy Estimators/ A. Kaltchenko, N. Timofeeva // International Journal of Information and Coding Theory. - 2009. - pp.39–56.

Оригинал-макет подготовлен  
в редакционно-издательском отделе ЯрГУ

Отпечатано на ризографе

Ярославский государственный университет  
150000 Ярославль, ул. Советская, 14.