

# Выбор метрики для оценки энтропии по методу ближайшей точки\*

**Е. А. ТИМОФЕЕВ**

Ярославский государственный университет  
им. П. Г. Демидова  
e-mail: TimofeevEA@gmail.com

УДК 519.72+519.234

**Ключевые слова:** энтропия, стохастический процесс, стационарность, непараметрическая оценка, метрика, смещение.

## Аннотация

Рассматривается задача повышения эффективности непараметрических оценок энтропии стационарного эргодического процесса. Изучаемый подход основан на методе ближайшей точки. Предлагается новый широкий класс метрик на пространстве правосторонних последовательностей, состоящих из символов конечного алфавита. Эти метрики имеют параметр — невозрастающую функцию. Показано, что при некоторых условиях предлагаемые оценки имеют малую дисперсию и что при специальном выборе параметров метрики можно уменьшить смещение.

## Abstract

*E. A. Timofeev, Selection of a metric for the nearest neighbor entropy estimators, Fundamentalnaya i prikladnaya matematika, vol. 18 (2013), no. 2, pp. 209–227.*

We consider the problem of improving the efficiency of the nonparametric entropy estimation for a stationary ergodic process. Our approach is based on the nearest-neighbor distances. We propose a broad class of metrics on the space of right-sided infinite sequences drawn from a finite alphabet. The new metric has a parameter which is a non-increasing function. We prove that, under certain conditions, our estimators have a small variance and show that a special selection of the metric parameters reduces the estimator's bias.

## 1. Введение

Работа посвящена повышению эффективности (точности) оценивания энтропии информационного источника с конечным числом состояний. Энтропия (энтропия на символ) является важнейшей характеристикой информационного источника. Поэтому задача оценивания энтропии представляет значительный интерес.

---

\*Работа выполнена при поддержке гранта Правительства РФ по постановлению № 220, договор 11.G34.31.0053.

Поскольку статистические характеристики информационного источника обычно не известны, наиболее широко используются так называемые *непараметрические* оценки энтропии. Однако аналитически определить точность этих оценок очень трудно, и известно очень мало результатов в этом направлении. В большинстве опубликованных работ по непараметрическому оцениванию энтропии доказывается только асимптотическая сходимость и приводятся результаты компьютерного тестирования.

Напомним, что для заданной выборки размера  $n$  наиболее важной характеристикой оценки  $h_n$  является её эффективность (точность), или квадратичная ошибка  $E(h_n - h)^2$ , где  $h$  — энтропия источника, а  $n$  — число измерений. Напомним соотношение

$$E(h_n - h)^2 = \text{Var } h_n + (Eh_n - h)^2,$$

где величина  $E(h_n - h)$  называется *смещением*.

Цель настоящей работы — построение новой оценки, основывающейся на методе ближайшей точки и его модификациях, у которой точность будет равна  $O(n^{-c})$  для широкого класса мер, где  $c > 0$  — некоторая константа.

Основная идея предлагаемой конструкции состоит в следующем.

Оценки энтропии, основанные на методе ближайшей точки, используют некоторую метрику. Будет введён довольно широкий класс так называемых *слабых* [3] метрик, для которых неравенство треугольника выполняется с некоторой константой  $C > 1$ .

Эта метрика применяется в построении оценки энтропии по методу ближайшей точки. Построенная оценка будет линейно зависеть от параметров метрики. А эти параметры подбираются так, чтобы минимизировать среднее отклонение.

Отметим, что задачи выбора параметров с линейными ограничениями довольно часто возникают в статистике [13].

Подчеркнём, что наиболее известные непараметрические оценки энтропии основаны либо на методе сжатия Лемпеля—Зива, либо на методе ближайшей точки. Для этих оценок доказана сходимость почти всюду (наиболее общие результаты приводятся, например, в [6, 9]). К сожалению, из-за медленной скорости сходимости (не более чем  $O(1/\ln n)$ ) их точность не является достаточной для большинства практических приложений, поскольку логарифм размера выборки ( $\log_2 n$ ) довольно небольшой:  $\log_2 n \leq 30 - 40$ . Этим объясняется поиск новых оценок с большей скоростью сходимости.

Следует отметить, что оценки энтропии, основанные на методе Лемпеля—Зива, очень трудны для теоретического исследования их сходимости и получения аналитических выражений для точности. Так, например, первоначальной мотивацией работы [1] была задача получения асимптотического значения скорости сходимости оценок энтропии, введённых Зивом [16], но авторы только смогли указать, что нахождение смещения и дисперсии очень трудно даже для простейшей меры Бернулли. На сегодняшний день нет опубликованных работ, в которых были бы найдены смещение или дисперсия оценки энтропии, построенной по методу Лемпеля—Зива.

Далее будут рассматриваться только оценки, построенные по методу ближайшей точки.

Заметим, что в то время как многие оценки применяются к одной длинной последовательности, оценки, основанные на методе ближайшей точки, применяются к большому числу независимых коротких последовательностей. Эти два подхода эквивалентны, если информационный источник имеет память размера порядка  $\ln n$ . Действительно, пусть дана выборка  $x_1, x_2, \dots, x_N$ ,  $x_i \in \mathcal{A}$ , где  $\mathcal{A}$  — алфавит информационного источника. Выберем достаточно большие числа  $m, l \approx C \ln n$ , положим  $n = \lceil (N - m)/(m + l) \rceil + 1$  и определим точки  $\xi_1, \dots, \xi_n$ , положив

$$\xi_i = (x_{(m+l)(i-1)+1}, x_{(m+l)(i-1)+2}, \dots, x_{(m+l)(i-1)+m}), \quad i = 1, 2, \dots, n.$$

Теперь кратко опишем основные известные результаты. Сразу отметим, что удобнее оценивать не энтропию  $h$ , а обратную величину  $1/h$ .

Будем считать, что информационный источник, или стационарный процесс, является инвариантной относительно сдвига эргодической мерой  $\mu$  на пространстве  $\Omega = A^{\mathbb{N}}$ , состоящем из правосторонних бесконечных последовательностей символов из конечного алфавита  $A$ , где  $\mathbb{N} = \{1, 2, \dots\}$ . Таким образом, бесконечную случайную последовательность будем рассматривать как случайную точку в  $\Omega$  распределённую по мере  $\mu$ .

Пусть

$$\rho_0(\mathbf{x}, \mathbf{y}) = \max\{e^{-k} : x_k \neq y_k\} - \quad (1)$$

метрика на пространстве  $\Omega$ , где  $\mathbf{x} = (x_1, x_2, \dots)$ ,  $\mathbf{y} = (y_1, y_2, \dots)$ .

Пусть

$$\xi_0 = (\xi_{01}, \dots, \xi_{0m}, \dots), \dots, \xi_n = (\xi_{n1}, \dots, \xi_{nm}, \dots) -$$

$n + 1$  случайных точек в  $\Omega$ . Тогда первая оценка (величины, обратной к энтропии), построенная по методу ближайшей точки [5], определяется следующим образом:

$$\tilde{h}_n = -\frac{1}{(n+1) \ln n} \sum_{j=0}^n \ln \left( \min_{i: i \neq j} \rho_0(\xi_i, \mathbf{x}_j) \right), \quad (2)$$

здесь и далее в работе все логарифмы будут рассматриваться с натуральным основанием.

Позже П. Шилдс [12] показал, что оценка Грассбергера не является состоятельной для общих эргодических процессов, но справедлива для неприводимых цепей Маркова.

Две модификации оценки Грассбергера были предложены в [14]. В обозначениях настоящей работы их можно записать как  $r_n^{(k)}(\rho) / \ln n$  (см. (27)) и  $\eta_n^{(k)}(\rho)$  (см. (26)), где  $\rho$  — некоторая метрика.

Для оценки  $r_n^{(k)}(\rho) / \ln n$   $L^1$ -сходимость и оценка дисперсии  $O(n^{-c})$  доказаны в [14] при довольно общих ограничениях на метрику  $\rho$  и меру  $\mu$ . В [6] для метрики (1) ограничения на меру существенно упрощены (см. (3)) и доказана

сходимость почти всюду. Также в [6] доказана оценка дисперсии  $O(n^{-c})$  для любого  $c < 1$ .

Для оценки  $\eta_n^{(k)}(\rho)$   $L^1$ -сходимость показана в [14] при более сильных ограничениях на меру и метрику.

Для метрики (1) вычислительный эксперимент, приведённый в [6], показал, что оценка (26) намного эффективнее, чем оценка  $r_n^{(k)}(\rho_0)/\ln n$ . Однако в следующей работе [7] для симметричных мер Бернулли было показано, что смещение оценки (26) является периодической функцией с периодом пропорциональным  $\ln n$ . В вычислительных экспериментах это смещение очень трудно заметить, поскольку его амплитуда не более чем  $10^{-6}$  для мер с небольшой энтропией ( $h < 3$ ).

В [15] для смещения найдено аналитическое выражение для марковских мер и метрики (1). Это смещение равно нулю, если логарифмы вероятностей перехода рационально несоизмеримы. В противном случае смещение является периодической функцией с периодом пропорциональным  $\ln n$ .

Этот результат показывает новую трудность в аналитических вычислениях смещения, а именно что смещение оценки может быть разрывной функцией от параметров меры.

Статья организована следующим образом:

- в разделе 2 приводится постановка задачи;
- в разделе 3 вводятся новые слабые метрики;
- в разделе 4 описаны оценки энтропии, основанные на введённых метриках;
- в разделе 5 изучаются статистики, используемые в методе ближайшей точки, и показывается, что дисперсия этих статистик оценивается как  $O(n^{-1} \ln^4 n)$  для довольно широкого класса мер и рассматриваемых слабых метрик;
- в разделе 6 для симметричной меры Бернулли показано, что существуют такие параметры метрик, при которых оценка является несмещённой;
- в разделе 7 приводится описание алгоритмов для быстрого нахождения оценок и обоснование их трудоёмкости.

## 2. Постановка задачи

Пусть  $\Omega = A^{\mathbb{N}}$  — пространство правосторонних бесконечных последовательностей символов из конечного алфавита  $A$  и  $\mu$  — инвариантная относительно сдвига эргодическая вероятностная мера на  $\Omega$ .

Пусть  $\xi_0, \xi_1, \dots, \xi_n$  — независимые случайные точки в  $\Omega$ , одинаково распределённые по мере  $\mu$ .

Требуется оценить энтропию меры  $\mu$ .

Дополнительно добавим следующее ограничение на меру:

найдутся  $a, b > 0$ , такие что  $\mu(C_n(\mathbf{x})) \leq be^{-an}$   
для любого  $n > 0$  для почти всех  $\mathbf{x} \in \Omega$ , (3)

где через

$$C_s(\mathbf{x}) = \{\mathbf{y} \in \Omega: y_1 = x_1, \dots, y_s = x_s\}$$

будут обозначаться цилиндры в пространстве  $\Omega$ .

Пусть случайная точка  $\xi = (\xi_1, \xi_2, \dots)$  в  $\Omega$  распределена по мере  $\mu$ . Напомним, что энтропия  $h$  меры  $\mu$  определяется следующим образом:

$$h = - \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E} \ln \mu(C_n(\xi)). \quad (4)$$

### 3. Метрики на пространстве последовательностей

В этом разделе будет описан новый широкий класс так называемых *слабых* метрик [3], для которых неравенство треугольника выполняется с некоторой константой  $C > 1$ .

Пусть  $\mathcal{A} = \{0, 1, \dots, A - 1\}$ , и будем считать, что  $A$  чётно.

Для точки  $\mathbf{x} = (x_1, x_2, \dots)$  из  $\Omega = \mathcal{A}^{\mathbb{N}}$ , через  $a\mathbf{x}$  будем обозначать точку  $(a, x_1, x_2, \dots)$ ,  $a \in \mathcal{A}$ .

Класс метрик  $\rho$  определяется следующим образом:

$$\rho(\mathbf{x}, \mathbf{y}) = e^{-\alpha(\mathbf{x}, \mathbf{y})}, \quad (5)$$

где функция  $\alpha(\mathbf{x}, \mathbf{y})$  определяется по следующим правилам:

$$\alpha(a\mathbf{x}, b\mathbf{y}) = \begin{cases} \alpha(\mathbf{x}, \mathbf{y}) + 1, & a = b, \\ \lambda_{a,b}(\alpha(\mathbf{x}, \mathbf{y})), & a \neq b. \end{cases} \quad (6)$$

Вспомогательные функции  $\lambda_{a,b}(t)$  неубывающие и удовлетворяют условию

$$0 \leq \lambda_{a,b}(t) = \lambda_{b,a}(t) \leq 1.$$

Более того, семейство функций  $\lambda_{a,b}(t)$  выбирается так, чтобы множество функций

$$S_a = \{\lambda_{a,0}(t), \dots, \lambda_{a,a-1}(t), \lambda_{a,a+1}(t), \dots, \lambda_{a,A-1}(t)\}$$

не зависело от  $a$  (т. е.  $S_0 = S_1 = \dots = S_{A-1}$ ) и  $|S_a| = A - 1$ . Другими словами, функции  $\lambda_{a,b}(t)$  можно интерпретировать как раскраску рёбер полного графа с  $A$  вершинами  $A - 1$  цветами. Отметим, что такая раскраска существует только для чётных  $A$ .

Выберем хорошо известную раскраску

$$\begin{aligned} \sigma_{0,1} = \sigma_{2,A-1} = \sigma_{3,A-2} = \dots = \sigma_{A/2,A/2+1} = 1, \\ \dots \\ \sigma_{0,k} = \sigma_{k-1,k+1} = \dots = \sigma_{1,2k-1} = \sigma_{2k,A-1} = \dots = \sigma_{A/2+k-1,A/2+k} = k, \quad (7) \\ \dots \\ \sigma_{0,A-1} = \sigma_{1,A-2} = \sigma_{2,A-3} = \dots = \sigma_{A/2-1,A/2} = A - 1 \end{aligned}$$

и положим

$$\lambda_{i,j}(t) = \lambda_{\sigma_{i,j}}(t), \quad 0 \leq i \neq j \leq A - 1, \quad (8)$$

где

$$\lambda_k(t) = \phi\left(\frac{k}{A} + \frac{1}{A}\phi^{-1}(t)\right), \quad k = 1, 2, \dots, A-1, \quad (9)$$

а  $\phi(t)$  — невозрастающая функция на полуинтервале  $(0, 1]$ , такая что

$$\phi(1) = 0, \quad \phi\left(\frac{t}{A}\right) = \phi(t) + 1. \quad (10)$$

Положим

$$\lambda_0(t) = t + 1, \quad (11)$$

тогда, очевидно, (9) выполняется и при  $k = 0$ .

Подчеркнём, что функции  $\lambda_k(t)$  такие, что

$$\lambda_0(t_0) > \lambda_1(t_1) > \dots > \lambda_{A-1}(t_{A-1}) \quad \text{для любых } t_0, t_1, \dots, t_{A-1} \geq 0. \quad (12)$$

Итак, можно переписать (6) как

$$\alpha(a\mathbf{x}, b\mathbf{y}) = \lambda_{\sigma_{a,b}}(\alpha(\mathbf{x}, \mathbf{y})),$$

где

$$\sigma_{a,a} = 0.$$

Применяя (9), получаем

$$\alpha(a\mathbf{x}, b\mathbf{y}) = \phi\left(\frac{\sigma_{a,b}}{A} + \frac{1}{A}\phi^{-1}(\alpha(\mathbf{x}, \mathbf{y}))\right). \quad (13)$$

Отметим, что функция  $\phi(t)$  является произвольной на интервале  $(1/A, 1]$  (невозрастающей и удовлетворяющей условию  $0 \leq \phi(t) \leq 1$ ).

Если  $\phi(t) = 0$ ,  $1/A < t \leq 1$ , то получаем метрику  $\rho_0$  (см. (1)).

Если  $\phi(t) = 1$ ,  $1/A \leq t \leq 1$ , получаем метрику  $(1/A)\rho_0$ .

Отметим также, что метрика  $\rho$  с произвольной функцией  $\phi(t)$  билипшицево эквивалентна метрике  $\rho_0$ , т. е.

$$e^{-1}\rho_0(\mathbf{x}, \mathbf{y}) \leq \rho(\mathbf{x}, \mathbf{y}) \leq \rho_0(\mathbf{x}, \mathbf{y}). \quad (14)$$

Следовательно,  $\rho$  является слабой метрикой [3], т. е. неравенство треугольника выполняется только с некоторой константой.

Каждая точка  $\mathbf{x}$  имеет бесконечное число координат, однако в практических вычислениях нужно использовать только конечное число из них. Сделаем это введением понятия сечения метрики, которое будет использовать только первые  $m$  координат каждой точки.

Определим  $\rho^{(m)}$  — сечение метрики  $\rho$ , положив

$$\rho^{(m)}(\mathbf{x}, \mathbf{y}) = e^{-\alpha^{(m)}(\mathbf{x}, \mathbf{y})}, \quad (15)$$

где

$$\begin{aligned} \alpha^{(0)}(\mathbf{x}, \mathbf{y}) &= 0, \\ \alpha^{(m)}(a\mathbf{x}, b\mathbf{y}) &= \lambda_{\sigma_{a,b}}(\alpha^{(m-1)}(\mathbf{x}, \mathbf{y})). \end{aligned} \quad (16)$$

**Утверждение 1.** Множество значений  $\alpha^{(m)}(\mathbf{x}, \mathbf{y})$  есть

$$\Phi_m = \left\{ \phi \left( \frac{N}{A^m} \right), N = 1, 2, \dots, A^m \right\}. \quad (17)$$

**Доказательство.** Проведём доказательство индукцией по  $m$ .

При  $m = 0$  утверждение очевидно.

Применяя определение (16), получаем

$$\Phi_m = \bigcup_{k=0}^{A-1} \lambda_k(\Phi_{m-1}). \quad (18)$$

По предположению индукции и (9), (11) имеем

$$\bigcup_{k=0}^{A-1} \lambda_k(\Phi_{m-1}) = \left\{ \phi \left( \frac{k}{A} + \frac{N}{A^m} \right), N = 1, 2, \dots, A^{m-1}, k = 0, 1, 2, \dots, A-1 \right\}.$$

Подставляя в (18), получаем (17).  $\square$

Для прикладных вычислений также нужно ограничить число параметров, которые используются для нахождения функции  $\phi(t)$ . Для этого введём сечение функции  $\phi(t)$ .

Определим  $\phi_l(t)$ , сечение функции  $\phi(t)$ , следующим образом:

$$\begin{aligned} \phi_l(t) &= \phi \left( \frac{N}{A^l} \right), \quad \frac{N-1}{A^l} < t \leq \frac{N}{A^l}, \quad N = A^{l-1} + 1, A^{l-1} + 2, \dots, A^l, \\ \phi_l \left( \frac{t}{A} \right) &= \phi_l(t) + 1, \quad 0 < t \leq 1. \end{aligned} \quad (19)$$

Следует подчеркнуть, что функции  $\phi_l(t)$  имеют множество значений

$$\left\{ \phi \left( \frac{N}{A^l} \right) + k, N = A^{l-1} + 1, \dots, A^l, k = 0, 1, \dots \right\}. \quad (20)$$

Следовательно, обратная функция  $\phi_l^{-1}(x)$  определена на этом множестве и имеет множество значений

$$\left\{ \frac{N}{A^{l+k}}, N = A^{l-1} + 1, \dots, A^l, k = 0, 1, \dots \right\}.$$

Применяя утверждение 1, получаем следствие.

**Следствие 1.** Пусть  $\phi(t) = \phi_l(t)$  и  $m \geq l$ . Тогда множество значений  $\alpha^{(m)}(\mathbf{x}, \mathbf{y})$  совпадает с

$$\begin{aligned} &\left\{ \phi \left( \frac{N}{A^l} \right) + k, N = A^{l-1} + 1, \dots, A^l, k = 0, 1, \dots, m-l \right\} \cup \\ &\cup \bigcup_{k=1}^{l-1} \left\{ \phi \left( \frac{N}{A^k} \right) + m-k, N = A^{k-1} + 1, \dots, A^k \right\} \cup \{m\}. \end{aligned} \quad (21)$$

**Утверждение 2.** Пусть  $\phi(t) = \phi_l(t)$ ,  $m \geq l$  и  $\alpha(\mathbf{x}, \mathbf{y}) < m - l + 1$ . Тогда

$$\alpha^{(m)}(\mathbf{x}, \mathbf{y}) = \alpha(\mathbf{x}, \mathbf{y}).$$

**Доказательство.** Итерируя (13), получаем

$$\alpha(\mathbf{x}, \mathbf{y}) = \phi_l \left( \frac{N}{A^j} + \frac{1}{A^j} \phi_l^{-1}(\alpha(\mathfrak{S}^j \mathbf{x}, \mathfrak{S}^j \mathbf{y})) \right), \quad (22)$$

где через  $\mathfrak{S}$  обозначается сдвиг на пространстве  $\Omega$ , а  $N$  — некоторое целое число,  $0 \leq N < A^j$ ,  $j \geq 1$ .

Поскольку  $\alpha(\mathbf{x}, \mathbf{y}) < m - l + 1$ , то для некоторого  $k \leq m - l$

$$\alpha(\mathbf{x}, \mathbf{y}) = \phi_l \left( \frac{N'}{A^l} \right) + k,$$

где  $A^{l-1} + 1 \leq N' \leq A^l$ . Подставляя  $j = l + k$  в (22), имеем

$$\alpha(\mathbf{x}, \mathbf{y}) = \phi_l \left( \frac{N}{A^{l+k}} + \frac{1}{A^{l+k}} \phi_l^{-1}(t) \right), \quad (23)$$

где  $t = \alpha(\mathfrak{S}^{l+k} \mathbf{x}, \mathfrak{S}^{l+k} \mathbf{y})$ . Поэтому

$$\phi_l \left( \frac{N'}{A^l} \right) + k = \phi_l \left( \frac{N}{A^{l+k}} + \frac{1}{A^{l+k}} \phi_l^{-1}(t) \right).$$

Отсюда следует, что  $N < A^l$  и

$$\phi_l \left( \frac{N'}{A^l} \right) = \phi_l \left( \frac{N}{A^l} + \frac{1}{A^l} \phi_l^{-1}(t) \right).$$

Поскольку  $0 < \phi_l^{-1}(t) \leq 1$ , имеем

$$N = N' - 1.$$

Итерируя (16), получаем

$$\alpha^{(m)}(\mathbf{x}, \mathbf{y}) = \phi_l \left( \frac{N}{A^{l+k}} + \frac{1}{A^{l+k}} \phi_l^{-1}(t_m) \right),$$

где  $t_m = \alpha^{(m-l-k)}(\mathfrak{S}^{l+k} \mathbf{x}, \mathfrak{S}^{l+k} \mathbf{y})$  и  $N$  такое же, как в (23). Следовательно,

$$\alpha^{(m)}(\mathbf{x}, \mathbf{y}) = \phi_l \left( \frac{N}{A^l} + \frac{1}{A^l} \phi_l^{-1}(t_m) \right) + k = \phi_l \left( \frac{N+1}{A^l} \right) + k = \alpha(\mathbf{x}, \mathbf{y}). \quad \square$$

#### 4. Оценки энтропии, построенные по методу ближайшей точки

Пусть  $n + 1$  точек  $\xi_0, \dots, \xi_n$  заданы своими первыми  $m$  координатами. Применим сечение  $\rho^{(m)}$  метрики (5)–(10).

Будем считать, что сечение функции (10) задаётся функцией  $\phi_l(t)$ ,  $l < m$ .



Пусть функция  $\phi_l(t)$  задана параметрами

$$\beta_i = \phi \left( \frac{A^{l-1} + i}{A^l} \right), \quad i = 1, 2, \dots, A^l - A^{l-1} - 1. \quad (24)$$

Подчеркнём, что

$$1 \geq \beta_1 > \beta_2 > \dots > \beta_{A^l - A^{l-1} - 1} \geq 0. \quad (25)$$

Кроме точек и метрики, будем использовать вспомогательный параметр  $k$ , который служит для контроля применимости. Оценки, полученные для различных значений  $k$ , являются оценками одной и той же величины.

Оценка  $\eta_n^{(k)}(\rho^{(m)})$  величины, обратной к энтропии  $1/h$ , определяется следующим образом [14]:

$$\eta_n^{(k)}(\rho^{(m)}) = k \left( r_n^{(k)}(\rho^{(m)}) - r_n^{(k+1)}(\rho^{(m)}) \right), \quad (26)$$

где

$$r_n^{(k)}(\rho^{(m)}) = \frac{1}{n+1} \sum_{j=0}^n \max_{i: i \neq j}^{(k)} \alpha^{(m)}(\xi_i, \xi_j) \quad (27)$$

и  $\max^{(k)}\{X_1, \dots, X_N\} = X_k$ , если  $X_1 \geq X_2 \geq \dots \geq X_N$ .

Применяя следствие 1, получаем основное свойство нашей метрики.

**Следствие 2.** Оценка  $\eta_n^{(k)}(\rho^{(m)})$  и статистика  $r_n^{(k)}(\rho^{(m)})$  являются линейными функциями от параметров  $\beta_i$ ,  $i = 1, 2, \dots, A^l - A^{l-1} - 1$ .

## 5. Статистические свойства оценок

В этом разделе будут изучены статистические свойства статистик (27) и оценок (26).

Обозначим открытый шар радиуса  $r$  с центром в точке  $\mathbf{x} \in \Omega$  через

$$B(\mathbf{x}, r, \rho) = \{\mathbf{y} \in \Omega: \rho(\mathbf{x}, \mathbf{y}) < r\}.$$

Обозначим обратную функцию к  $t = \mu(B(\mathbf{x}, e^{-u}, \rho))$  через  $u = \nu(t, \mathbf{x}, \rho)$ , т. е.

$$\nu(t, \mathbf{x}, \rho) = \inf\{u: \mu(B(\mathbf{x}, e^{-u}, \rho)) < t\}. \quad (28)$$

Для краткости будем использовать обозначения

$$\nu(t, \mathbf{x}) = \nu(t, \mathbf{x}, \rho), \quad B(\mathbf{x}, r) = B(\mathbf{x}, r, \rho).$$

Следует подчеркнуть, что функции  $t = \mu(B(\mathbf{x}, r))$  и  $r = \nu(t, \mathbf{x})$  могут иметь интервалы постоянства и точки разрыва (интервал постоянства одной функции соответствует разрыву другой).

**Лемма 1.** Пусть  $\mu$  — инвариантная относительно сдвига эргодическая мера на  $\Omega$  и  $\rho$  — метрика (5). Тогда для  $\mu$ -почти всех точек  $\mathbf{x} \in \Omega$

$$\lim_{r \rightarrow 0} \frac{\ln \mu(B(\mathbf{x}, r, \rho))}{\ln r} = h,$$

где  $h$  — энтропия  $\mu$ .

**Доказательство.** Сначала рассмотрим случай  $\rho = \rho_0$ . Шары в метрике  $\rho_0$  являются цилиндрами, т. е.

$$B(\mathbf{x}, r, \rho_0) = C_n(\mathbf{x}), \quad e^{-n-1} < r \leq e^{-n}.$$

Поэтому

$$\lim_{r \rightarrow 0} \frac{\ln \mu(B(\mathbf{x}, r, \rho_0))}{\ln r} = - \lim_{n \rightarrow \infty} \frac{\ln \mu(C_n(\mathbf{x}))}{n}.$$

Применяя теорему Шеннона—Макмиллана—Бреймана [10, 2.10], получаем, что для  $\mu$ -почти всех точек  $\mathbf{x} \in \Omega$

$$\lim_{r \rightarrow 0} \frac{\ln \mu(B(\mathbf{x}, r, \rho_0))}{\ln r} = - \lim_{n \rightarrow \infty} \frac{\ln \mu(C_n(\mathbf{x}))}{n} = h.$$

Теперь рассмотрим произвольную метрику (5). Поскольку  $\rho$  билипшицево эквивалентна (14) метрике (1), имеем

$$B(\mathbf{x}, e^{-1}r, \rho_0) \subset B(\mathbf{x}, r, \rho) \subset B(\mathbf{x}, er, \rho_0).$$

Поэтому

$$\mu(B(\mathbf{x}, e^{-1}r, \rho_0)) \leq \mu(B(\mathbf{x}, r, \rho)) \leq \mu(B(\mathbf{x}, er, \rho_0)).$$

Итак, для  $\mu$ -почти всех точек  $\mathbf{x} \in \Omega$

$$\lim_{r \rightarrow 0} \frac{\log \mu(B(\mathbf{x}, r, \rho))}{\log r} = h. \quad \square$$

**Лемма 2.**

$$Er_n^{(k)}(\rho) = n \binom{n-1}{k-1} \int_{\Omega} \int_0^1 \nu(u, \mathbf{x}, \rho) u^{k-1} (1-u)^{n-k} du d\mu(\mathbf{x}). \quad (29)$$

**Доказательство.** Поскольку  $P\{\alpha(\mathbf{x}, \xi_i) > t\} = \mu(B(\mathbf{x}, e^{-t}, \rho))$ , имеем

$$\begin{aligned} E \max_{1 \leq j \leq n}^{(k)} \alpha(\mathbf{x}, \xi_j) &= n \binom{n-1}{k-1} \int_0^{\infty} t \mu(B(\mathbf{x}, e^{-t}, \rho))^{k-1} \times \\ &\quad \times \left(1 - \mu(B(\mathbf{x}, e^{-t}, \rho))\right)^{n-k} d_t \mu(B(\mathbf{x}, e^{-t}, \rho)). \end{aligned}$$

Следовательно,

$$\begin{aligned} Er_n^{(k)}(\rho) &= E \max_{1 \leq j \leq n}^{(k)} \alpha(\xi_0, \xi_j) = E[E \max_{1 \leq j \leq n}^{(k)} \alpha(\xi_0, \xi_j) \mid \xi_0] = \\ &= n \binom{n-1}{k-1} \int_{\Omega} \int_0^{\infty} t \mu(B(\mathbf{x}, e^{-t}, \rho))^{k-1} \times \\ &\quad \times \left(1 - \mu(B(\mathbf{x}, e^{-t}, \rho))\right)^{n-k} d_t \mu(B(\mathbf{x}, e^{-t}, \rho)) d\mu(\mathbf{x}). \end{aligned}$$

Применяя формулу замены переменных, получаем (29). □

**Лемма 3.** Пусть мера  $\mu$  удовлетворяет условию (3). Пусть  $\rho$  — метрика (5) с  $\phi(t) = \phi_l(t)$ . Тогда существуют такие константы  $c_1, c_2$ , что верно неравенство

$$0 < \text{Er}_n^{(k)}(\rho) - \text{Er}_n^{(k)}(\rho^{(m)}) \leq c_1 n^{-1} \quad \text{для } m \geq l + c_2 \log n. \quad (30)$$

**Доказательство.** Согласно утверждению 2

$$B(\mathbf{x}, e^{-t}, \rho) = B(\mathbf{x}, e^{-t}, \rho^{(m)}), \quad t \leq m - l,$$

поэтому

$$\nu(u, \mathbf{x}, \rho) = \nu(u, \mathbf{x}, \rho^{(m)}), \quad u \geq \varepsilon_m,$$

где

$$\varepsilon_m = \mu(B(\mathbf{x}, e^{-m+l}, \rho)).$$

Применив лемму 2, получим

$$\begin{aligned} \text{Er}_n^{(k)}(\rho) - \text{Er}_n^{(k)}(\rho^{(m)}) &= \\ &= n \binom{n-1}{k-1} \int_{\Omega} \int_0^{\varepsilon_m} [\nu(u, \mathbf{x}, \rho) - \nu(u, \mathbf{x}, \rho^{(m)})] u^{k-1} (1-u)^{n-k} du d\mu(\mathbf{x}) \leq \\ &\leq n \binom{n-1}{k-1} \int_{\Omega} \int_0^{\varepsilon_m} \nu(u, \mathbf{x}, \rho) u^{k-1} du d\mu(\mathbf{x}). \end{aligned}$$

Используя условие (3) и неравенство (14), имеем

$$\mu(B(\mathbf{x}, e^{-t}, \rho)) \leq c_3 e^{-at}.$$

Поэтому

$$\nu(u, \mathbf{x}, \rho) \leq -\frac{1}{a} \ln \frac{u}{c_3}, \quad \varepsilon_m \leq c_3 e^{-a(m-l)}.$$

Подставляя эти неравенства, получаем

$$\text{Er}_n^{(k)}(\rho) - \text{Er}_n^{(k)}(\rho^{(m)}) \leq -\frac{n^k}{a} \int_0^{\varepsilon_m} \ln \frac{u}{c_3} u^{k-1} du.$$

Вычислив интеграл, получим

$$\text{Er}_n^{(k)}(\rho) - \text{Er}_n^{(k)}(\rho^{(m)}) = O\left(n^k m e^{-ak(m-l)}\right).$$

Возьмём  $c_2 > (k+1)/(ak)$ . Тогда неравенство (30) следует из полученного равенства.  $\square$

Применяя лемму 1 и повторяя доказательство теоремы 1 из [14], получаем следующее утверждение.

**Утверждение 3.** Пусть  $\xi_0, \dots, \xi_n - n + 1$  независимых точек в пространстве  $\Omega$ , распределённых по мере  $\mu$  и  $k = O(\ln n)$ . Тогда

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

**Следствие 3.** Пусть выполнены условия леммы 3 и  $k = O(\ln n)$ . Тогда

$$\lim_{n \rightarrow \infty} \frac{\mathbb{E} r_n^{(k)}(\rho^{(m)})}{\ln n} = \frac{1}{h}.$$

**Теорема 4.** Пусть  $r_n^{(k)}(\rho^{(m)})$  — статистика, определённая в (27). Тогда выполняется следующее неравенство:

$$\text{Var } r_n^{(k)}(\rho^{(m)}) \leq \frac{m^2(km+1)^2 A^2}{4(n+1)}. \quad (31)$$

**Доказательство.** Применим метод Макдиармида [11].

Введём функцию

$$f: \Omega^{n+1} \rightarrow \mathbb{R},$$

положив

$$f(\mathbf{x}_0, \dots, \mathbf{x}_n) = \frac{1}{n+1} \sum_{j=0}^n \max_{i: i \neq j}^{(k)} \alpha^{(m)}(\mathbf{x}_i, \mathbf{x}_j).$$

Для того чтобы применить метод Макдиармида, нужно показать, что  $f$  удовлетворяет неравенству

$$\sup_{\mathbf{x}_0, \dots, \mathbf{x}_n, \mathbf{y} \in \Omega} |f(\mathbf{x}_0, \dots, \mathbf{x}_n) - f(\mathbf{x}_0, \dots, \mathbf{x}_{i-1}, \mathbf{y}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n)| \leq c_i \quad (32)$$

для  $0 \leq i \leq n$ . Докажем это неравенство для

$$c_i = \frac{m(km+1)A}{n+1}, \quad 0 \leq i \leq n. \quad (33)$$

Для краткости введём следующие обозначения:

$$\begin{aligned} \mathbf{X} &= (\mathbf{x}_0, \dots, \mathbf{x}_n), \\ \tilde{\mathbf{X}} &= (\mathbf{x}_0, \dots, \mathbf{x}_{i-1}, \mathbf{y}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n), \\ g_j(\mathbf{X}) &= \max_{l: l \neq j}^{(k)} \alpha^{(m)}(\mathbf{x}_l, \mathbf{x}_j). \end{aligned}$$

Положим

$$J = \{j \neq i: g_j(\mathbf{X}) \neq g_j(\tilde{\mathbf{X}})\}.$$

Поскольку

$$\begin{aligned} g_j(\mathbf{X}) &= g_j(\tilde{\mathbf{X}}), \quad j \notin J, \quad j \neq i, \\ |g_j(\mathbf{X}) - g_j(\tilde{\mathbf{X}})| &\leq m \quad \text{для всех } j, \end{aligned}$$

имеем

$$|f(\mathbf{X}) - f(\tilde{\mathbf{X}})| \leq \frac{m(|J|+1)}{n+1}.$$

По утверждению 5  $|J| \leq kmA$ , поэтому выполняется (33).

Как показано в [2], метод Макдиармида даёт следующую оценку дисперсии  $f(\xi_0, \dots, \xi_n)$ :

$$\text{Var}[f(\xi_0, \dots, \xi_n)] \leq \frac{1}{4} \sum_{i=0}^n c_i^2. \quad (34)$$

Подставляя (33) в (34), получаем (31). □

**Следствие 4.**

$$P\{|r_n^{(k)}(\rho^{(m)}) - Er_n^{(k)}(\rho^{(m)})| > \delta\} \leq 2e^{-2(n+1)\delta^2/m^2(km+1)^2A^2}. \quad (35)$$

**Доказательство.** Это неравенство получается применением следующего неравенства из [11]:

$$P\{|f(\xi_0, \dots, \xi_n) - Ef(\xi_0, \dots, \xi_n)| > \delta\} \leq 2e^{-2\delta^2/\sum_{i=0}^n c_i^2},$$

где  $c_i$  заданы в (33). Подставляя (33) в это неравенство, получаем

$$P\{|f(\xi_0, \dots, \xi_n) - Ef(\xi_0, \dots, \xi_n)| > \delta\} \leq 2e^{-2(n+1)\delta^2/m^2(km+1)^2A^2}.$$

Подставляя  $r_n^{(k)}(\rho^{(m)}) = f(\xi_0, \dots, \xi_n)$ , получаем (35). □

**Следствие 5.** Пусть мера  $\mu$  удовлетворяет условию (3),  $\rho$  — метрика (5) с  $\phi(t) = \phi_l(t)$ ,  $m \geq c_1 \ln n + l$ ,  $m \leq c_2 \ln n$ , где  $c_1, c_2$  — некоторые константы, и  $k = O(\ln n)$ . Тогда  $Er_n^{(k)}(\rho^{(m)})/\ln n$  почти всюду сходится к  $1/h$ .

Применяя теорему 4 и неравенство  $\text{Var}(X + Y) \leq (\sqrt{\text{Var } X} + \sqrt{\text{Var } Y})^2$ , получаем следующее утверждение.

**Утверждение 4.** Пусть  $\rho$  — метрика (5) с  $\phi(t) = \phi_l(t)$ ,  $m \geq c_1 \ln n + l$ ,  $m \leq c_2 \ln n$ , где  $c_1, c_2$  — некоторые константы,  $k, A$  фиксированные. Тогда

$$\text{Var} \eta_m^{(k)}(\rho^{(m)}) = O\left(\frac{\ln^4 n}{n}\right).$$

## 6. Симметричная мера Бернулли

Покажем, что для симметричной меры Бернулли (с  $A$  равновероятными символами) можно найти такую функцию  $\phi(t)$  для метрики (5)–(10), что  $\eta_m^{(k)}(\rho)$  является несмещённой оценкой.

Для симметричной меры Бернулли функция  $\nu(t, \omega)$  не зависит от точки  $\omega$ . Поэтому  $\chi(t) = \nu(t, \omega)$  задаётся следующим рекуррентным уравнением:

$$\chi(t) = \begin{cases} \chi(At) + 1, & t \leq \frac{1}{A}, \\ \phi\left(\frac{k}{A} + \frac{1}{A}\phi^{-1}(\chi(At - k))\right), & \frac{k}{A} < t \leq \frac{k+1}{A}, \quad k = 1, 2, \dots, A-1. \end{cases} \quad (36)$$

Нетрудно убедиться, что решением этого уравнения является  $\chi(t) = \phi(t)$ .

Покажем, что для

$$\chi(t) = \phi(t) = -\frac{\ln t}{\ln A}$$

$\eta_n^{(k)}(\rho)$  будет несмещённой оценкой.

Подставляя  $\chi(t)$  в (29), имеем

$$Er_n^{(k)}(\rho) = -\frac{n!}{(k-1)!(n-k)!\ln A} \int_0^1 \ln(t) t^{k-1} (1-t)^{n-k} dt.$$

Применяя [4, 4.253.1], получаем

$$Er_n^{(k)}(\rho) = \frac{\psi(n+1) - \psi(k)}{\ln A}, \quad (37)$$

где

$$\psi(t) = \frac{d}{dt} \ln \Gamma(t) -$$

пси-функция. Следовательно,

$$E\eta_n^{(k)}(\rho) = k \frac{\psi(k+1) - \psi(k)}{\ln A} = \frac{1}{\ln A}. \quad (38)$$

## 7. Алгоритмы

В этом разделе будет описано несколько алгоритмов для нахождения статистики (27).

### 7.1. Поиск $k$ -ближайшей

Пусть  $\xi_1, \dots, \xi_n$  —  $n$  слов  $\xi_i = \xi_{i1} \dots \xi_{im}$ ,  $\xi_{ij} \in \mathcal{A}$ ,  $i = 1, 2, \dots, n$ .

Для данного слова  $\xi_0$  и целого  $k$  нужно найти следующий максимум:

$$\max_{1 \leq i \leq n}^{(k)} \alpha^{(m)}(\xi_i, \xi_0).$$

Все слова  $\xi_1, \dots, \xi_n$  будем хранить как бор  $T$  [8]. Бор — это  $A$ -арное дерево, вершинами которого служат векторы длины  $A$ . Будем обозначать через  $T$  массив размера  $A \times N$ , столбцами которого служат эти  $N$  векторов. Таким образом,  $T(i, j)$  будет вершиной на следующем уровне от вершины  $i$ , если существует слово в словаре с заданным префиксом и следующим символом  $j$ , иначе  $T(i, j) = 0$ . Припишем каждой вершине  $i$  бора  $T$  дополнительный параметр  $L(i)$ , равный числу листьев в подборе с корнем  $i$ .

Отметим, что в нашем боре  $T$  все листья находятся на уровне  $m$ . Будем считать, что номер каждого листа равняется индексу  $i$  соответствующего слова  $\xi_i$ .

Метрика (5) задаётся функцией  $\phi(t)$  (10) и перестановками (7)

$$\sigma_i = (\sigma_{i,0}, \dots, \sigma_{i,A-1}).$$

Будем также использовать обратные перестановки

$$\bar{\sigma}_i = (\bar{\sigma}_{i,0}, \dots, \bar{\sigma}_{i,A-1}), \quad i = 0, 1, \dots, A-1.$$

Приведём псевдокод алгоритма поиска.

- Полагаем  $V := \text{root}$ .
- **for**  $l := 1$  **to**  $m-1$  **do**
- **for**  $i := 0$  **to**  $A-1$  **do**
- **if**  $L(T(\bar{\sigma}(\xi_{0,l}, i), V)) \geq k$
- **then**  $V := T(\bar{\sigma}(\xi_{0,l}, i), V)$ ; **break**;
- **else**  $k := k - L(T(\bar{\sigma}(\xi_{0,l}, i), V))$ .

Номер вершины  $V$  (на уровне  $m$ ) является индексом слова  $\xi_V$ , для которого

$$\begin{aligned} \alpha^{(m)}(\xi_V, \xi_0) &= \\ &= \phi \left( \frac{1}{A} \sigma_{\xi_{V,1}, \xi_{0,1}} + \frac{1}{A} \left( \frac{1}{A} \sigma_{\xi_{V,2}, \xi_{0,2}} + \dots + \frac{1}{A} \left( \frac{1}{A} \sigma_{\xi_{V,m}, \xi_{0,m}} + \frac{1}{A} \right) \dots \right) \right). \end{aligned} \quad (39)$$

Корректность алгоритма вытекает из следующего простого свойства метрики.

**Утверждение 5.** Пусть  $Y$  — такое слово, что

$$\alpha^{(m)}(Y, \xi_0) = \max_{1 \leq i \leq n}^{(k)} \alpha^{(m)}(\xi_i, \xi_0),$$

пусть  $Y$  лежит в подборе  $T_0$ , таком что число листьев  $T_0$  не менее  $k$ , и пусть  $r_0, \dots, r_{A-1}$  — потомки корня  $T_0$ , занумерованные индексами соответствующих функций  $\lambda_i(t)$ . Тогда  $Y$  принадлежит подбору с корнем  $r_i$ , таким что

$$L(r_0) + \dots + L(r_{i-1}) < k, \quad L(r_0) + \dots + L(r_i) \geq k.$$

Доказательство следует непосредственно из свойства (12).

Временная трудоёмкость алгоритма равна  $\mathcal{O}(mA)$ . Требуемая память —  $\mathcal{O}(mAn)$ .

## 7.2. Эффективный алгоритм нахождения $r_n^{(k)}(\rho^{(m)})$

Для того чтобы найти  $r_n^{(k)}(\rho^{(m)})$ , построим бор  $T$  с дополнительным параметром  $L$  из слов  $\xi_0, \dots, \xi_n$ . Затем для  $i = 0, 1, \dots, n$  найдём  $(k+1)$ -ближайшее слово в боре  $T$  для слова  $\xi_i$ . Суммируя, получаем  $r_n^{(k)}(\rho^{(m)})$ .

Временная сложность равна  $\mathcal{O}(mnA)$ .

## 7.3. Эффективный алгоритм нахождения $r_n^{(k)}(\rho^{(m)})$ для всех $n$

В этом разделе рассмотрим следующую задачу. Пусть  $\xi_0, \dots, \xi_n$  —  $n+1$  слов  $\xi_i = \xi_{i1} \dots \xi_{im}$ ,  $\xi_{ij} \in \mathcal{A}$ ,  $i = 0, 1, \dots, n$ . Требуется найти значения  $r_j^{(i)}(\rho^{(m)})$  для всех  $i, j$ , где  $k \leq j \leq n$ ,  $1 \leq i \leq k$ .

Введём вспомогательный параметр

$$R(i, j) = \max_{1 \leq s \leq n}^{(i)} \alpha^{(m)}(\xi_s, \xi_j). \quad (40)$$

Для данной таблицы слов  $\xi$  строим вышеописанный бор  $T$  и находим вспомогательный параметр  $R$ .

Приведём псевдокод алгоритма.

- Полагаем  $R := 0$ ,  $T := \emptyset$ .
- **for**  $j := 1$  **to**  $n$  **do**
  - вставляем слово  $\xi_j$  в бор  $T$ ;
  - находим  $R(i, j)$  для  $i = 1, 2, \dots, k$ , применяя алгоритм из раздела 7.1;
  - пусть  $p_0, \dots, p_m$  — путь в боре от листа  $j$  к корню и  $Q := 0$ ;
  - for**  $t := 1$  **to**  $m$  **do**
    - for** каждого сына  $S$  вершины  $p_t$  **do**
    - if**  $L(S) \leq k$  **then**
    - for** каждого листа  $V$  of  $S$  **do**
    - for**  $i := L(S)$  **to**  $k$  **do**
      - полагаем  $Q(i) := Q(i) - R(i, V)$ ;
      - пересчитываем  $R(i, V)$ , применяя алгоритм из раздела 7.1;
      - полагаем  $Q(i) := Q(i) + R(i, V)$ ;
  - **for**  $i := 1$  **to**  $k$  **do**
    - полагаем  $r_j^{(i)}(\rho^{(m)}) = (1/j)((j-1)r_{j-1}^{(i)}(\rho^{(m)}) + Q(i) + R(i, j))$ .

Корректность алгоритма вытекает из следующего утверждения.

**Утверждение 6.** Пусть  $V$  — такая вершина в боре  $T$ , что  $L(V) > k$ , и  $T_0$  — подбор с корнем  $V$ . Тогда для любого слова  $\xi_i$ , соответствующего листу  $i$  подбора  $T_0$ , имеем

$$\max_{0 \leq j \neq i \leq n}^{(k)} \alpha^{(m)}(\xi_i, \xi_j) = \alpha^{(m)}(\xi_i, \xi_s),$$

где лист  $s$  также принадлежит подбору  $T_0$ .

Доказательство непосредственно следует из утверждения 5.

Временная трудоёмкость равна  $\mathcal{O}(mk^2An)$ .

#### 7.4. Символическое вычисление $r_n^{(k)}(\rho^{(m)})$ для всех $n$

Предположим, что сечение  $\phi_l(t)$  задано параметрами  $\beta_i$  (24), где  $\beta_i$  неизвестны,  $i = 1, 2, \dots, M$ ,

$$M = A^l - A^{l-1} - 1. \quad (41)$$

Алгоритм из раздела 7.1 применим и в этом случае, но в (39) получим

$$\alpha^{(m)}(\xi_V, \xi_0) = U\beta_l + V,$$

где  $U, V$  не зависят от  $\beta$ .



Подставляя в (27), имеем

$$r_n^{(k)}(\rho^{(m)}) = \frac{1}{n}D_{n,k} + \frac{1}{n} \sum_{s=1}^M C_{n,k,s} \beta_s. \quad (42)$$

В этом подразделе рассмотрим задачу нахождения величин  $D_{j,i}$ ,  $C_{j,i,s}$  для всех  $i, j$ , где  $k \leq j \leq n$ ,  $1 \leq i \leq k$ ,  $1 \leq s \leq M$ .

Введём вспомогательные параметры  $I, U, V$ :

$$\max_{1 \leq s \leq n}^{(i)} \alpha^{(m)}(\xi_s, \xi_j) = U(i, j) \beta_{I(i,j)} + V(i, j). \quad (43)$$

Пусть бор  $T$  задан вышеописанным образом. Тогда по алгоритму последовательно находим вспомогательные параметры  $I, U, V$  и  $D, C$ .

Приведём псевдокод алгоритма.

- Полагаем  $C := 0$ ,  $D := 0$ ,  $U := 0$ ,  $V := 0$ .
- **for**  $j := 1$  **to**  $n$  **do**
- **for**  $i := 1$  **to**  $k$  **do**
  - находим  $I(i, j)$ ,  $U(i, j)$ ,  $V(i, j)$ , применяя алгоритм из раздела 7.1;
  - полагаем  $\tilde{D}(i) := V(i, j)$ ;
  - полагаем  $\tilde{C}(i, I(i, j)) := U(i, j)$ ;
- пусть  $p_0, \dots, p_m$  — путь в боре  $T$  от листа  $j$  до корня;
- for**  $t := 1$  **to**  $m$  **do**
  - for** каждого сына  $S$  вершины  $p_t$  **do**
  - if**  $L(S) \leq k$  **then**
  - for** каждого листа  $F$  вершины  $S$  **do**
  - for**  $i := L(S)$  **to**  $k$  **do**
    - полагаем  $\tilde{D}(i) := \tilde{D}(i) - V(i, F)$ ,
    - $\tilde{C}(i, I(i, F)) := \tilde{C}(i, I(i, F)) - U(i, F)$ ;
    - пересчитываем  $I(i, F)$ ,  $U(i, F)$ ,  $V(i, F)$ ,
    - применяя алгоритм из раздела 7.1;
    - полагаем  $\tilde{D}(i) := \tilde{D}(i) + V(i, F)$ ,
    - $\tilde{C}(i, I(i, F)) := \tilde{C}(i, I(i, F)) + U(i, F)$ ;
- **for**  $i := 1$  **to**  $k$  **do**
  - полагаем  $D(j, i) := D(j, i) + \tilde{D}(i)$ ;
  - for**  $s := 1$  **to**  $M$  **do**
    - полагаем  $C(j, i, s) := C(j, i, s) + \tilde{C}(i, s)$ .

Корректность алгоритма следует из утверждения 6.

Временная трудоёмкость равна  $\mathcal{O}(mk^2An + kMn)$ .

### 7.5. Нахождение параметров метрики

Параметры метрики будем выбирать так, чтобы минимизировать среднее отклонение. По следствию 2 оценка является линейной функцией от параметров, поэтому получаем задачу минимизации положительной квадратичной формы на симплексе (25). Действительно, по (42) имеем

$$\eta_n^{(k)}(\rho^{(m)}) = \eta_n^{(k)}(\rho_0^{(m)}) + \sum_{i=1}^M \beta_i \tilde{C}_{n,k,i}, \quad (44)$$

где

$$\tilde{C}_{n,k,i} = k(C_{n,k,i} - C_{n,k+1,i}).$$

Будем минимизировать функцию

$$F(\beta) = \frac{1}{n-k+1} \sum_{j=k}^n \left( \eta_j^{(k)}(\rho^{(m)}) - \bar{\eta}_j^{(k)}(\rho^{(m)}) \right)^2, \quad (45)$$

где

$$\bar{\eta}_n^{(k)}(\rho^{(m)}) = \frac{1}{n-k+1} \sum_{j=k}^n \eta_j^{(k)}(\rho^{(m)}). \quad (46)$$

Подставляя (44) в (45), (46), получаем, что  $F(\beta)$  является квадратичной формой от своих параметров. Следовательно, задача минимизации этой функции является задачей минимизации положительной квадратичной формы на симплексе (25).

### 7.6. Алгоритм нахождения оценки энтропии

Алгоритм состоит из двух этапов.

1. Выбираем часть заданных строк.  
Находим коэффициенты линейной функции (44).  
Находим минимум положительной квадратичной формы (45) на симплексе (25).
2. Находим оценку (26) энтропии по оставшимся заданным строкам с параметрами, найденными на первом этапе.

Автор пользуется возможностью выразить признательность В. Л. Дольникову за полезные обсуждения.

## Литература

- [1] Aldous D., Shields P. A diffusion limit for a class of randomly-growing binary trees // Probab. Th. Rel. Fields. — 1988. — Vol. 79. — P. 509—542.

- [2] Devroye L. Exponential inequalities in nonparametric estimation // *Nonparametric Functional Estimation and Related Topics* / G. Roussas, ed. — Dordrecht: Kluwer Academic, 1991. — P. 31–44.
- [3] Deza M., Deza T. *Encyclopedia of Distances*. — Berlin: Springer, 2009.
- [4] Gradshteyn I. S., Ryzhik I. M. *Table of Integrals, Series, and Products*. — London: Academic Press, 1994.
- [5] Grassberger P. Estimating the information content of symbol sequences and efficient codes // *IEEE Trans. Inform. Theory*. — 1989. — Vol. 35. — P. 669–675.
- [6] Kaltchenko A., Timofeeva N. Entropy estimators with almost sure convergence and an  $O(n^{-1})$  variance // *Adv. Math. Commun.* — 2008. — Vol. 2. — P. 1–13.
- [7] Kaltchenko A., Timofeeva N. Rate of convergence of the nearest neighbor entropy estimator // *AEU — Int. J. Electron. Commun.* — 2010. — Vol. 64. — P. 75–79.
- [8] Knuth D. E. *The Art of Computer Programming*. Vol. 3. *Sorting and Searching*. — Reading: Addison-Wesley, 1975.
- [9] Kontoyiannis I., Suhov Yu. M. Prefixes and the entropy rate for long-range sources // *Probability Statistics and Optimization* / F. P. Kelly, ed. — New York: Wiley, 1994. — P. 89–98.
- [10] Martin N., England J. *Mathematical Theory of Entropy*. — Cambridge: Cambridge Univ. Press, 1984.
- [11] McDiarmid C. On the method of bounded differences // *Surveys in Combinatorics*. — Cambridge: Cambridge Univ. Press, 1989. — P. 148–188.
- [12] Shields P. C. Entropy and prefixes // *Ann. Probab.* — 1992. — Vol. 20. — P. 403–409.
- [13] Silvapulle M. J., Sen P. K. *Constrained Statistical Inference: Inequality, Order and Shape Restrictions*. — New York: Wiley, 2005.
- [14] Timofeev E. A. Statistical estimation of measure invariants // *St. Petersburg Math. J.* — 2006. — Vol. 17, no. 3. — P. 527–551.
- [15] Timofeev E. A. Bias of a nonparametric entropy estimator for Markov measures // *J. Math. Sci.* — 2011. — Vol. 176, no. 2. — P. 255–269.
- [16] Ziv J., Lempel A. Compression of individual sequences by variable rate coding // *IEEE Trans. Inform. Theory*. — 1978. — Vol. 24. — P. 530–536.

