

ПРИКЛАДНАЯ ТЕОРИЯ КОДИРОВАНИЯ

УДК 519.728

ТЕОРЕТИЧЕСКИ ЭФФЕКТИВНОЕ АСИМПТОТИЧЕСКИ ОПТИМАЛЬНОЕ УНИВЕРСАЛЬНОЕ КОДИРОВАНИЕ ЧАСТИЧНО ОПРЕДЕЛЁННЫХ ИСТОЧНИКОВ

Л. А. Шоломов

ФИЦ «Информатика и управление» РАН, г. Москва, Россия

Частично определённый источник порождает независимо с некоторыми вероятностями символы заданного основного алфавита и неопределённый символ. Кодирование источника должно обеспечить точное воспроизведение основных символов, а неопределённые символы допускают замену (доопределение) любыми основными символами. Кодирование считается эффективным, если имеет полиномиальную оценку сложности кодирования и декодирования. Оно асимптотически оптимально, если обеспечивает среднюю длину кода, асимптотически равную энтропии источника. Кодирование универсально, если оно не зависит от вероятностей символов источника. Описан метод эффективного асимптотически оптимального универсального кодирования частично определённых источников.

Ключевые слова: *недоопределённый источник, частично определённый источник, универсальное кодирование, полиномиальный метод, энтропия источника, квазиэнтропия слова, частотный класс, комбинаторная энтропия, представительное множество.*

DOI 10.17223/20710410/47/4

THEORETICALLY EFFECTIVE ASYMPTOTICALLY OPTIMAL UNIVERSAL CODING OF PARTIALLY DEFINED SOURCES

L. A. Sholomov

Federal Research Center “Computer Science and Control” of RAS, Moscow, Russia

E-mail: levshol@mail.ru

Let $A_0 = \{a_i : i \in M\}$ be a finite alphabet of basic symbols, $A = \{a_T : T \in \mathcal{T}\}$, $\mathcal{T} \subseteq 2^M$, — an alphabet of underdetermined symbols. Any symbol a_i , $i \in T$, is considered to be a specification of the symbol a_T . The symbol a_M , denoted by $*$, is called indefinite. An underdetermined source X generates symbols $a_T \in A$ independently with probabilities p_T . The entropy of the source X is the quantity

$$\mathcal{H}(X) = \min_Q \left(- \sum_{T \in \mathcal{T}} p_T \log \sum_{i \in T} q_i \right),$$

where the minimum is taken over the set of probability vectors $Q = (q_i, i \in M)$, $\log x = \log_2 x$. For source X , we consider separable block coding with a block length of n . Encoding K of underdetermined words $v \in A^n$ should ensure that the code $K(v)$

of word v allow one to recover some specification of v . Coding is universal if it is independent of the probabilities p_T , $T \in \mathcal{T}$. Characteristic of coding quality is average code length

$$\bar{l}^{(n)} = \frac{1}{n} \sum_{v \in A^n} p(v) |K(v)|,$$

where $p(v) = p_{T_1} \dots p_{T_n}$ is the probability of the appearance of the word $v = a_{T_1} \dots a_{T_n}$ at the source output, $|K(v)|$ is the codeword length for v .

Earlier, the author established that for arbitrary source X and for any block coding method, the inequality $\bar{l}^{(n)} \geq \mathcal{H}(X)$ holds, and that there is universal block coding, for which $\bar{l}^{(n)} \leq \mathcal{H}(X) + O\left(\frac{\log n}{n}\right)$. The upper bound here is obtained by random coding which only establishes existence of the estimation, but does not provide an appropriate algorithm. Important for applications are issues of complexity of procedures. Coding method considered theoretically effective if the complexities of coding and decoding are estimated by a some polynomial of the size n of problem. This paper presents a polynomial coding method for underdetermined sources whose alphabet has the form $A = A_0 \cup \{*\} = \{a_0, a_1, \dots, a_{m-1}, *\}$. Such sources are called partially defined. For them, entropy can be explicitly expressed:

$$\mathcal{H}(P) = (1 - p_*) \log(1 - p_*) - \sum_{i=0}^{m-1} p_i \log p_i.$$

The main content of the paper is the proof of the following result.

For a partially defined source X , there is an universal method of block coding that provides an estimate of the average length of the code

$$\bar{l}^{(n)} \leq \mathcal{H}(P) + O\left(\frac{\log \log n}{\log^{1/2} n}\right)$$

and allows you to encode and decode using RAM-programs with complexity $O(n^2)$. An analogue of this result is also valid for other types of undetermined sources whose entropy is explicitly representable.

Keywords: *underdetermined source, partially defined source, universal coding, polynomial method, source entropy, quasientropy of a word, frequency class, combinatorial entropy, representative set.*

Введение

Под недоопределёнными данными понимаются последовательности недоопределённых символов. Каждому такому символу соответствует некоторое множество символов основного алфавита, любым из которых он может быть замещен (доопределён). Считается, что последовательности порождаются источником, генерирующим символы недоопределённого алфавита A независимо с некоторыми вероятностями. Задача кодирования такого источника состоит в том, чтобы каждой порождаемой им последовательности сопоставить двоичный код, позволяющий восстановить некоторое доопределение последовательности (но не её саму). Рассматривается блочное кодирование источника, при котором последовательности разбиваются на блоки некоторой длины n и производится кодирование блоков — слов длины n в алфавите A . Качество кодирования характеризуется средней по блокам длиной кода, приходящейся на символ источника [1]. Кодирование универсально, если оно не зависит от вероятностей символов.

С недоопределённым источником X связывается энтропия $\mathcal{H}(X)$ [2], роль которой подобна роли энтропии Шеннона для полностью определённых источников. В [3] (см. также [2]) доказано, что при любом способе кодирования недоопределённого источника X средняя длина кода не меньше $\mathcal{H}(X)$ и что существует его универсальное кодирование, обеспечивающее среднюю длину кода, асимптотически равную $\mathcal{H}(X)$. Этот результат получен с использованием метода случайного кодирования, который лишь устанавливает факт существования указанной асимптотически оптимальной оценки средней длины кода, но не обеспечивает соответствующего алгоритма и потому не позволяет утверждать что-либо нетривиальное о сложности кодирования.

Вопросы сложности процедур важны для приложений. В теории сложности принято считать процедуру эффективной, если она является детерминированной и её сложность (число элементарных операций) оценивается сверху полиномом от размера задачи. Это понятие эффективности фактически не зависит от модели вычислений, ибо для основных моделей переход от одной из них к другой связан с не более чем полиномиальным изменением числа операций. Применительно к задачам кодирования теоретическая эффективность подразумевает, что полиномиальными должны быть процедуры кодирования и декодирования. Чтобы иметь возможность говорить о степени полинома в оценке сложности, следует уточнить модель вычислений. В качестве модели будем использовать РАМ-программу [4] с подходящим набором операций. Отметим, что оценки сложности носят асимптотический характер и теоретическая эффективность ещё не означает, что процедуры являются «практически хорошими» при реальных значениях параметров.

В данной работе представлен полиномиальный (квадратичный) метод асимптотически оптимального кодирования недоопределённых источников, имеющих более частный вид. Они порождают лишь полностью определённые и неопределённые символы и названы частично определёнными. Отметим, что в приложениях чаще всего встречаются недоопределённые данные такого вида. В случае двоичного основного алфавита понятия недоопределённого и частично определённого источников совпадают.

Распространение описанного в работе эффективного метода кодирования частично определённых источников на случай недоопределённых источников общего вида затруднено тем, что в общем случае энтропия $\mathcal{H}(X)$ задана неявно как минимум некоторой функции и её явное выражение удаётся найти лишь в редких случаях, к которым относится и случай частично определённых источников. Явное выражение энтропии используется при кодировании. Метод применим и к другим типам источников с явно выражимой энтропией (например, к рассмотренным в [5]). Отметим, что в предлагаемом эффективном методе остаточный член в асимптотической оценке средней длины кода значительно хуже, чем в неэффективном методе из [3].

Приведём известные результаты, относящиеся к проблематике данной работы. Имеется тесная связь рассмотренной в работе вероятностной задачи кодирования недоопределённого источника с детерминированной задачей кодирования класса слов с заданными частотами вхождения недоопределённых символов (см. лемму 10 данной работы). В терминах [6] можно детерминированную задачу кодирования недоопределённого слова понимать как задачу построения двоичной программы вычисления некоторого его доопределения. В случае двоичных недоопределённых слов в качестве способа вычисления довольно естественно рассматривать схему из логических элементов (неветвящуюся программу [4, 7]), реализующую булеву функцию, последовательность значений которой совпадает с этим словом. Двоичное представление этой схемы можно считать кодом слова, и длина этого кода определяется сложностью схемы.

Идея воспринимать схему как код реализуемой функции восходит к К. Шеннону [8] и лежит в основе развитого им мощностного метода получения нижних оценок сложности схем. О. Б. Лупанов [9] разработал общий метод синтеза схем — принцип локального кодирования, позволивший получить асимптотически точные оценки сложности для многих классов функций и систем. В частности, этим методом Лупанов реализовал асимптотически наилучшим образом полностью определённые функции от n аргументов, принимающие N_1 единичных значений (для всех N_1 , удовлетворяющих минимальному естественному ограничению). Построение асимптотически наилучших схем этот принцип связывает с кодированием функций, удовлетворяющим ряду требований. Некоторые приёмы такого кодирования использованы в данной работе.

Работа Э. И. Нечипорука [10] об асимптотически наилучшей реализации недоопределённых двоичных матриц с заданным числом булевых элементов вентильными схемами глубины 2 явилась пионерской в области синтеза недоопределённых управляющих систем. Метод этой работы с использованием принципа локального кодирования позволяет строить асимптотически наилучшие схемы из логических элементов для булевых функций от n аргументов, определённых на N наборах при больших значениях $N = N(n)$. В [11] предложены методы повышения степени определённости функций путём сведения к функциям от меньшего числа аргументов, определённым примерно на том же числе наборов. Это позволило решить задачу наилучшего синтеза схем в широком диапазоне значений параметра N . Окончательный результат получен в [12], где использован другой способ повышения степени определённости функций.

Детерминированной задаче кодирования двоичных частично определённых слов с заданными частотами появления символов соответствует задача схемной реализации частично определённой функции, заданной на N наборах и принимающей N_0 нулевых и N_1 единичных значений. Полное решение задачи построения асимптотически наилучших реализаций при всех значениях параметров (N, N_0, N_1) , удовлетворяющих минимальному ограничению, получено А. В. Чашкиным [13]. Результаты для некоторых диапазонов значений (N, N_0, N_1) были опубликованы ранее в [14–16]. Вклад каждой из этих работ в решение общей задачи охарактеризован в [13]. К этой проблематике примыкает результат из [17] об асимптотически минимальной реализации вентильными схемами глубины 2 многозначных частично определённых матриц с заданными частотами символов.

В работе [18] рассмотрены функции, значениями которых являются недоопределённые символы общего вида, и описан асимптотически наилучший метод их схемной реализации (при некотором двоичном представлении символов). Другое решение той же задачи — более структурированное — имеется в [19]. Рассматриваемая в этих работах постановка задачи связана с несколько упрощённым понятием энтропии, которое использует лишь мощности доопределяющих множеств для символов и не учитывает их структуру. Это понятие согласовано с задачей протыкания системы комбинаторных кубоидов [18], но в задаче кодирования недоопределённых слов и источников даёт завышенные значения.

В перечисленных работах решается в разных постановках задача построения асимптотически наилучших по сложности схем для недоопределённых функций, но не рассматриваются вопросы сложности самих процедур построения схем. Развитые там методы включают либо вероятностные, либо неэффективные комбинаторные построения и потому не могут служить основой эффективных методов кодирования недоопределённых слов. Исключение составляет работа [20], в которой описан эффективный способ получения двоичной программы для вычисления булевой функции, заданной

на N наборах. Эта программа имеет асимптотически минимальную длину и строится со сложностью, почти квадратично зависящей от размера исходных данных.

1. Основные понятия и формулировка результата

Задан конечный алфавит $A_0 = \{a_0, a_1, \dots, a_{m-1}\}$ *основных* символов. Каждому непустому подмножеству T множества $M = \{0, 1, \dots, m-1\}$ соответствует символ a_T , называемый *недоопределённым*. Основные символы a_i будем также рассматривать как недоопределённые символы, соответствующие одноэлементным подмножествам $\{i\} \subseteq M$. *Доопределением символа a_T* считается всякий основной символ a_i , $i \in T$. Символ a_M , доопределимый любым основным символом, называется *неопределённым* и обозначается $*$.

Выделена система $\mathcal{T} \subseteq 2^M$ некоторых непустых множеств $T \subseteq M$ и ей сопоставлен *недоопределённый алфавит* $A = \{a_T : T \in \mathcal{T}\}$. Под *доопределением слова* $v = a_{T_1} \dots a_{T_l} \in A^l$ понимается любое слово $w = a_{i_1} \dots a_{i_l} \in A_0^l$, полученное из v заменой каждого символа каким-либо его доопределением.

Рассматривается источник X , порождающий символы $a_T \in A$ независимо с вероятностями p_T , $\sum_{T \in \mathcal{T}} p_T = 1$. Он обозначается (A, P) , где $P = (p_T, T \in \mathcal{T})$, и называется *недоопределённым источником*. Задавшись некоторым набором вероятностей $Q = (q_i, i \in M)$ основных символов, введём функцию

$$\mathcal{H}(P, Q) = - \sum_{T \in \mathcal{T}} p_T \log \sum_{i \in T} q_i \quad (1)$$

(здесь и дальше логарифмы двоичные). *Энтропией источника X* назовём величину

$$\mathcal{H}(X) = \mathcal{H}(P) = \min_Q \mathcal{H}(P, Q). \quad (2)$$

Для недоопределённых источников она играет роль, подобную той, какую для полностью определённых источников играет энтропия Шеннона (подробнее см. в [2]).

Дальше считаем, что вероятности p_T всех символов $a_T \in A$ строго положительны, ибо в случае $p_T = 0$ символ a_T можно исключить из рассмотрения (а множество T — из \mathcal{T}). Очевидно, что $\mathcal{H}(X) \geq 0$. Нас интересует случай, когда энтропия $\mathcal{H}(X)$ строго положительна. С учётом соглашения $p_T > 0$ необходимое и достаточное условие положительности энтропии (имеющееся, например, в [2]) приобретает вид $\bigcap_{T \in \mathcal{T}} T = \emptyset$.

Рассматривается *блоковое кодирование* последовательностей, порождаемых недоопределённым источником X . Выходные последовательности источника разбиваются на куски некоторой длины n (*блоки*) и каждый блок — слово v длины n в алфавите A — заменяется его кодом $K(v)$ при некотором способе кодирования K . Считается, что кодирование является двоичным и *разделимым*. Последнее означает, что по произвольной конкатенации $K(v_1) \dots K(v_s)$ кодовых слов образующие её слова $K(v_1), \dots, K(v_s)$ восстанавливаются однозначно. К кодированию K недоопределённых слов $v \in A^n$ предъявляется требование, чтобы оно обеспечивало возможность нахождения по коду $K(v)$ какого-либо доопределения слова v . Кодирование источника (A, P) называется *универсальным* (при заданной длине n блоков), если оно не зависит от набора вероятностей P .

Качество кодирования K будем характеризовать *средней длиной кода*

$$\bar{l}_K^{(n)} = \frac{1}{n} \sum_{v \in A^n} p(v) |K(v)|, \quad (3)$$

приходящейся на символ источника [1]. Здесь $p(v) = p_{T_1} \dots p_{T_n}$ — вероятность порождения источником X слова $v = a_{T_1} \dots a_{T_n}$; $|K(v)|$ — длина кодового слова для v .

В [3] получен следующий результат (см. также [2]).

Теорема 1. Для произвольного недоопределённого источника X с $\mathcal{H}(X) > 0$

1) при любом способе K блочного кодирования имеет место неравенство

$$\bar{l}_K^{(n)} \geq \mathcal{H}(X);$$

2) существует универсальное блочное кодирование K , для которого

$$\bar{l}_K^{(n)} \leq \mathcal{H}(X) + O\left(\frac{\log n}{n}\right).$$

Верхняя оценка теоремы 1 получена методом случайного кодирования, который лишь устанавливает факт существования указанной в теореме асимптотически оптимальной оценки средней длины кода, но не обеспечивает соответствующего алгоритма. В данной работе представлен полиномиальный метод кодирования недоопределённых источников более частного вида. Недоопределённый источник $X = (A, P)$ называется *частично определённым*, если его алфавит имеет вид $A = A_0 \cup \{*\} = \{a_0, a_1, \dots, a_{m-1}, *\}$. Ему соответствует набор вероятностей $P = (p_0, p_1, \dots, p_{m-1}, p_*)$. Для частично определённого источника энтропия (2) допускает явное выражение (см. [2], а также п. 3° леммы 1 данной работы):

$$\mathcal{H}(P) = (1 - p_*) \log(1 - p_*) - \sum_{i=0}^{m-1} p_i \log p_i.$$

При условии положительности всех p_i эта величина строго положительна тогда и только тогда, когда $m \geq 2$.

Основное содержание работы составляет доказательство следующего результата.

Теорема 2. Для частично определённого источника $X = (A, P)$ с положительной энтропией

1) при любом способе K блочного кодирования справедливо неравенство

$$\bar{l}_K^{(n)} \geq \mathcal{H}(P);$$

2) имеется метод K универсального блочного кодирования, обеспечивающий оценку средней длины кода

$$\bar{l}_K^{(n)} \leq \mathcal{H}(P) + O\left(\frac{\log \log n}{\log^{\frac{1}{2}} n}\right)$$

и допускающий кодирование и декодирование РАМ-программами сложности $O(n^2)$.

Замечание 1. При доказательстве верхней оценки средней длины кода остаточный член фактически получен в виде $O\left(\phi(n)\left(\frac{\log \log n}{\log n}\right)^{1/2}\right)$, где $\phi(n)$ — произвольная растущая функция (см. лемму 11). Оценка теоремы соответствует значению $\phi(n) = (\log \log n)^{1/2}$. Её можно улучшать, беря в качестве $\phi(n)$ меньшую функцию.

Распространение теоремы 2 на общий случай недоопределённых источников затруднено тем, что доказательство использует явный вид энтропии $\mathcal{H}(X)$, который удаётся найти лишь в редких случаях.

2. Квазиэнтропия недоопределённых слов и ее свойства

Основной результат работы относится к частично определённым алфавитам (и источникам), но изложение будем вести, как правило, для недоопределённых алфавитов общего вида, учитывая возможность распространения результата работы на другие случаи. Если факт относится лишь к частично определённым данным, это будем отмечать явно.

Для слова $v \in A^n$ обозначим через $r_T(v)$ число появлений в нем символа a_T , $\sum_{T \in \mathcal{T}} r_T(v) = |v|$, и положим $\mathbf{r}(v) = (r_T(v), T \in \mathcal{T})$. Квазиэнтропией слова v назовём величину $h(v) = |v| \mathcal{H}\left(\frac{\mathbf{r}(v)}{|v|}\right)$. С учётом (1) и (2) она может быть переписана в виде

$$h(v) = \min_Q \left\{ - \sum_{T \in \mathcal{T}} r_T(v) \log \sum_{i \in \mathcal{T}} q_i \right\}. \quad (4)$$

Отметим, что точным аналогом принятого для полностью определённых данных понятия квазиэнтропии (см., например, [21]) является $\mathcal{H}\left(\frac{\mathbf{r}(v)}{|v|}\right)$. Но величина $|v| \mathcal{H}\left(\frac{\mathbf{r}(v)}{|v|}\right)$ оказывается более удобной, поскольку измеряется в тех же единицах, что и алгоритмическая энтропия (сложность по Колмогорову [6]), а также комбинаторная (см. далее).

Сначала приведём некоторые свойства энтропии $\mathcal{H}(X)$, которые понадобятся для изучения квазиэнтропии. Более детальное рассмотрение этих (и других) свойств энтропии имеется в [2].

Лемма 1.

1°) Энтропия $\mathcal{H}(P)$ неотрицательна и не превосходит $\log m$.

2°) Функция $\mathcal{H}(P)$ вогнута, т. е. для любых P, P' и числа $\theta, 0 \leq \theta \leq 1$, выполнено

$$\mathcal{H}(\theta P + (1 - \theta)P') \geq \theta \mathcal{H}(P) + (1 - \theta) \mathcal{H}(P').$$

3°) Энтропии источника $X = (A, P)$ и образованного из него исключением неопределённого символа $*$ $X' = (A', P')$, где $A' = A \setminus \{*\}$, $P' = (p'_T, T \in \mathcal{T} \setminus \{M\})$, $p'_T = \frac{p_T}{1 - p_*}$, связаны соотношением

$$\mathcal{H}(P) = (1 - p_*) \mathcal{H}(P').$$

4°) Для частично определённого источника (A, P) , где $A = \{a_0, a_1, \dots, a_{m-1}, *\}$, $P = (p_0, p_1, \dots, p_{m-1}, p_*)$, энтропия достигается в (2) на наборе

$$Q = \left(\frac{p_0}{1 - p_*}, \frac{p_1}{1 - p_*}, \dots, \frac{p_{m-1}}{1 - p_*} \right) \quad (5)$$

и принимает значение

$$\mathcal{H}(P) = (1 - p_*) \log(1 - p_*) - \sum_{i=0}^{m-1} p_i \log p_i. \quad (6)$$

Доказательство. Начнём с утверждений п. 1°. Неотрицательность энтропии очевидна. Верхнюю границу энтропии получим, вычислив $\mathcal{H}(P, Q)$ на наборе $Q = (1/m, \dots, 1/m)$:

$$\mathcal{H}(P) \leq - \sum_T p_T \log \frac{|T|}{m} = \log m - \sum_T p_T \log |T| \leq \log m.$$

Для доказательства п. 2° рассмотрим набор Q , на котором достигается минимум функции $\mathcal{H}(\theta P + (1 - \theta)P', Q)$. Имеем

$$\mathcal{H}(\theta P + (1 - \theta)P') = -\theta \sum_{T \in \mathcal{T}} p_T \log \sum_{i \in T} q_i - (1 - \theta) \sum_{T \in \mathcal{T}} p'_T \log \sum_{i \in T} q_i \geq \theta \mathcal{H}(P) + (1 - \theta) \mathcal{H}(P'),$$

что и требовалось доказать.

Докажем п. 3°. Из того, что для любого набора вероятностей $Q = (q_i, i \in M)$ выполнено $\log \sum_{i \in M} q_i = 0$, следует равенство

$$-\sum_{T \in \mathcal{T}} p_T \log \sum_{i \in T} q_i = -(1 - p_*) \sum_{T \in \mathcal{T} \setminus \{M\}} \frac{p_T}{1 - p_*} \log \sum_{i \in T} q_i.$$

Взяв минимум по Q , приходим к утверждению п. 3°.

Применим к частично определённому источнику из п. 4° утверждение п. 3°. Имеем

$$\mathcal{H}(P) = (1 - p_*) \mathcal{H}(P') = (1 - p_*) \min_Q \left\{ -\sum_{i \in M} \frac{p_i}{1 - p_*} \log q_i \right\}. \quad (7)$$

Как известно, указанный минимум достигается на наборе (5). Подставляя это значение Q в (7) и преобразуя с учётом равенства $\sum_{i \in M} p_i = 1 - p_*$, приходим к (6). ■

Следующая лемма 2 содержит необходимые для дальнейшего свойства квазиэнтропии слов.

Лемма 2.

- 1°) Квазиэнтропия $h(v)$ неотрицательна и не превосходит $n \log m$.
- 2°) Если слово v_2 получено из v_1 перестановкой символов, то $h(v_2) = h(v_1)$.
- 3°) Если $\mathbf{r}(v_2)$ отличается от $\mathbf{r}(v_1)$ лишь в компоненте $r_*(\cdot)$, то $h(v_2) = h(v_1)$.
- 4°) Квазиэнтропия конкатенации $v_1 v_2$ удовлетворяет неравенству

$$h(v_1 v_2) \geq h(v_1) + h(v_2).$$

- 5°) Если v_2 — продолжение слова v_1 , то $h(v_2) \geq h(v_1)$.
- 6°) Если va_T — результат приписывания к слову v символа a_T , то

$$h(va_T) \leq h(v) + \log(e|v|).$$

- 7°) Квазиэнтропия слова v в частично определённом алфавите $A = \{a_0, a_1, \dots, a_{m-1}, *\}$ задаётся выражением

$$h(v) = (|v| - r_*(v)) \log(|v| - r_*(v)) - \sum_{i \in M} r_i(v) \log r_i(v) \quad (8)$$

и достигается в (4) на наборе

$$Q = \left(\frac{r_0(v)}{|v| - r_*(v)}, \frac{r_1(v)}{|v| - r_*(v)}, \dots, \frac{r_{m-1}(v)}{|v| - r_*(v)} \right). \quad (9)$$

Доказательство. Пункт 1° следует из п. 1° леммы 1.

Пункт 2° вытекает из определения квазиэнтропии и того, что для слов v_2 и v_1 с одним и тем же составом символов имеет место $|v_2| = |v_1|$ и $\mathbf{r}(v_2) = \mathbf{r}(v_1)$.

Докажем п. 3°. Обозначим через v'_1 и v'_2 слова, полученные из v_1 и v_2 удалением всех символов $*$. По условиям п. 3° имеем $|v'_1| = |v_1| - r_*(v_1) = |v_2| - r_*(v_2) = |v'_2|$ и $\mathbf{r}(v'_1) = \mathbf{r}(v'_2)$, а потому

$$h(v'_1) = h(v'_2). \quad (10)$$

Из п. 4° леммы 1 при $p_T = \frac{r_T(v_1)}{|v_1|}$ с учётом равенства $1 - \frac{r_*(v_1)}{|v_1|} = \frac{|v'_1|}{|v_1|}$ выводим

$$\mathcal{H}\left(\frac{\mathbf{r}(v_1)}{|v_1|}\right) = \frac{|v'_1|}{|v_1|} \mathcal{H}\left(\frac{\mathbf{r}(v'_1)}{|v'_1|}\right).$$

Домножив обе части на $|v_1|$, приходим к равенству $h(v_1) = h(v'_1)$. Аналогично доказывается, что $h(v_2) = h(v'_2)$. Из этих равенств и (10) получаем $h(v_2) = h(v_1)$.

Для доказательства п. 4° воспользуемся вогнутостью энтропии (п. 3° леммы 1) при $P = \frac{\mathbf{r}(v_1)}{|v_1|}$, $P' = \frac{\mathbf{r}(v_2)}{|v_2|}$ и $\theta = \frac{|v_1|}{|v_1| + |v_2|}$.

$$\text{Имеем } \frac{|v_1|}{|v_1| + |v_2|} \mathcal{H}\left(\frac{\mathbf{r}(v_1)}{|v_1|}\right) + \frac{|v_2|}{|v_1| + |v_2|} \mathcal{H}\left(\frac{\mathbf{r}(v_2)}{|v_2|}\right) \leq \mathcal{H}\left(\frac{\mathbf{r}(v_1) + \mathbf{r}(v_2)}{|v_1| + |v_2|}\right) = \mathcal{H}\left(\frac{\mathbf{r}(v_1 v_2)}{|v_1 v_2|}\right).$$

Домножив обе части на $|v_1| + |v_2| = |v_1 v_2|$, получаем $h(v_1) + h(v_2) \leq h(v_1 v_2)$.

Чтобы доказать утверждение п. 5°, представим v_2 в виде $v_2 = v_1 w$, после чего в соответствии с п. 4° придём к $h(v_2) \geq h(v_1) + h(w)$ и воспользуемся неотрицательностью квазиэнтропии $h(w)$ (п. 1°).

Докажем п. 6°. Пусть $v = a_{T_1} \dots a_{T_{|v|}}$, $v' = va_T$. Обозначим через $Q = (q_0, q_1, \dots, q_{m-1})$ набор, на котором в (4) достигается квазиэнтропия $h(v)$. Возьмём некоторое значение $s \in T$ и образуем набор $Q' = (q'_0, q'_1, \dots, q'_{m-1})$, где $q'_s = \frac{(|v| - 1)q_s}{|v|} + \frac{1}{|v|}$, $q'_i = \frac{(|v| - 1)q_i}{|v|}$ для остальных $i \in M$. Набор Q' удовлетворяет условию

$$\sum_{i \in M} q'_i = \frac{|v| - 1}{|v|} \sum_{i \in M} q_i + \frac{1}{|v|} = \frac{|v| - 1}{|v|} + \frac{1}{|v|} = 1.$$

С учётом этого имеем

$$\begin{aligned} h(v') &\leq - \sum_{i \in M} r_{T_i}(v') \log \sum_{j \in T_i} q'_j = - \sum_{i \in M} r_{T_i}(v) \log \sum_{j \in T_i} q'_j - \log \sum_{j \in T} q'_j \leq \\ &\leq - \sum_{i \in M} r_{T_i}(v) \log \sum_{j \in T_i} \frac{|v| - 1}{|v|} q_j - \log \frac{1}{|v|} = - \sum_{i \in M} r_{T_i}(v) \log \sum_{j \in T_i} q_j - |v| \log \frac{|v| - 1}{|v|} + \log |v| = \\ &= h(v) - |v| \log \left(1 - \frac{1}{|v|}\right) + \log |v| \leq h(v) + \log e + \log |v| = h(v) + \log(e|v|). \end{aligned}$$

Утверждение п. 7° вытекает из п. 4° леммы 1 при $P = \frac{\mathbf{r}(v)}{|v|}$ с учётом равенства $\sum_{i \in M} r_i(v) = |v| - r_*(v)$. ■

Замечание 2. Выражение (8) может быть записано в виде, не зависящем от $|v|$:

$$h(v) = \left(\sum_{i \in M} r_i(v) \right) \log \sum_{i \in M} r_i(v) - \sum_{i \in M} r_i(v) \log r_i(v). \quad (11)$$

3. Частотные классы недоопределённых слов и их комбинаторная энтропия

Пусть задано конечное множество S недоопределённых слов в алфавите A . Скажем, что некоторое множество слов в алфавите A_0 доопределяет S , если в нём найдётся доопределение каждого слова из S . Обозначим через $N(S)$ минимальную мощность множества, доопределяющего S . Величину $\log N(S)$ будем называть комбинаторной энтропией множества слов S .

Для заданного набора $\mathbf{r} = (r_T, T \in \mathcal{T})$ натуральных чисел положим $l = \sum_{T \in \mathcal{T}} r_T$ и обозначим через $\mathcal{K}_l(\mathbf{r})$ класс всех слов длины l в алфавите A , в которых символ $a_T \in A$ встречается r_T раз (т.е. с частотой r_T/l). Такие классы называют частотными. Вместо записи $N(\mathcal{K}_l(\mathbf{r}))$ будем использовать $N_l(\mathbf{r})$. Комбинаторной энтропии класса $\mathcal{K}_l(\mathbf{r})$ соответствует обозначение $\log N_l(\mathbf{r})$. Всякое доопределяющее множество для класса $\mathcal{K}_l(\mathbf{r})$ будем называть доопределением класса $\mathcal{K}_l(\mathbf{r})$.

Все слова $v \in \mathcal{K}_l(\mathbf{r})$ имеют одинаковую квазиэнтропию $h(v) = l\mathcal{H}(\mathbf{r}/l)$, которую будем обозначать $h_l(\mathbf{r})$ и называть квазиэнтропией класса $\mathcal{K}_l(\mathbf{r})$.

Следующий результат из [3] (см. также [2, 22]) показывает, что квазиэнтропия частотного класса приближает его комбинаторную энтропию с точностью $O(\log l)$.

Утверждение 1. Справедливы оценки¹

$$h_l(\mathbf{r}) - C_1 \log l \leq \log N_l(\mathbf{r}) \leq h_l(\mathbf{r}) + C_2 \log l.$$

Верхняя оценка здесь получена методом случайного кодирования, который доказывает лишь существование доопределяющего множества для класса $\mathcal{K}_l(\mathbf{r})$, обеспечивающего эту оценку, но не даёт алгоритма. Ту же оценку комбинаторной энтропии мы получим ниже конструктивным построением доопределяющего множества с использованием градиентной процедуры (называемой также жадным алгоритмом, методом наискорейшего спуска, методом вычерпывания).

При осуществлении градиентной процедуры доопределения множества слов S используется некоторое базовое множество D всюду определённых слов, из которого заимствуются доопределения слов множества S . Паре (S, D) сопоставляется таблица с элементами 0 и 1, в которой строки соответствуют словам $v \in S$, столбцы — словам $w \in D$, и на пересечении строки v и столбца w помещается 1, если w доопределяет v , и 0 в противном случае. В терминах этой таблицы градиентная процедура доопределения множества S представляет собой последовательность шагов, на каждом из которых в таблице, полученной к данному шагу, выбирается столбец с максимальным числом единиц и вычёркиваются строки, содержащие в нём единицы. Совокупность выбранных столбцов к моменту, когда все строки окажутся вычеркнутыми, соответствует доопределяющему множеству для S . При оценке его мощности используется следующее утверждение, доказанное рядом авторов как обобщение леммы Нечипорука из [10] (оно имеется, например, в [23]).

Утверждение 2. Во всякой таблице с элементами 0 и 1, имеющей b строк, d столбцов и содержащей в каждой строке не менее a единиц, градиентная процедура позволяет выделить

$$k \leq \frac{d}{a} \left(\ln \frac{ba}{d} + 1 \right) + 1$$

столбцов так, что для любой строки имеется выделенный столбец, содержащий в ней 1.

¹Всюду в работе буквой C с индексом обозначается константа, абсолютная либо зависящая от мощности алфавитов A_0 и A .

Замечание 3. Градиентная процедура выделения столбцов полиномиальна относительно площади bd таблицы.

В градиентной процедуре построения доопределяющего множества для класса $\mathcal{K}_l(\mathbf{r})$ в качестве базового множества используется некоторый частотный класс $\mathcal{K}_l(\mathbf{m})$, $\mathbf{m} = (m_i, i \in M)$, слов в алфавите A_0 . Для нахождения набора \mathbf{m} параметров этого класса понадобится следующее свойство точки Q , на которой достигается энтропия в (2). Его доказательство имеется в [2, 3].

Утверждение 3. Точка $Q = (q_i, i \in M)$ минимума функции $\mathcal{H}(P, Q)$ удовлетворяет условиям

$$\sum_{T \in \mathcal{T}: i \in T} \frac{r_T q_i}{\sum_{j \in T} q_j} = q_i, \quad i \in M. \quad (12)$$

Докажем теперь основной результат данного раздела.

Лемма 3. Существует частотный класс $\mathcal{K}_l(\mathbf{m})$, при использовании которого в качестве базового множества градиентная процедура доопределения класса $\mathcal{K}_l(\mathbf{r})$ обеспечивает оценку

$$\log N_l(\mathbf{r}) \leq h_l(\mathbf{r}) + C_3 \log l. \quad (13)$$

Доказательство. Пусть $Q = (q_i, i \in M)$ — точка, в которой достигается квазиэнтропия $h_n(\mathbf{r})$, т. е. выполнено

$$-\sum_{T \in \mathcal{T}} r_T \log \sum_{i \in T} q_i = h_l(\mathbf{r}). \quad (14)$$

Введём величины s_{Ti} ($T \in \mathcal{T}, i \in M$), положив

$$s_{Ti} = \begin{cases} \frac{r_T q_i}{l \sum_{j \in T} q_j}, & i \in T, \\ 0, & i \notin T. \end{cases} \quad (15)$$

Величины s_{Ti} удовлетворяют условиям

$$\begin{aligned} \sum_i l s_{Ti} &= r_T \sum_{i \in T} \frac{q_i}{\sum_{j \in T} q_j} = r_T, \\ \sum_{T, i} l s_{Ti} &= \sum_T \sum_i l s_{Ti} = \sum_T r_T = l \end{aligned} \quad (16)$$

(если область изменения параметра T либо i при суммировании не указана, считается, что T пробегает множество \mathcal{T} , i — множество M . То же относится к параметрам величин, участвующих в произведении). Кроме того, поскольку на наборе Q , на котором достигается квазиэнтропия $h_l(\mathbf{r})$, достигается и энтропия $\mathcal{H}(\mathbf{r}/l)$, получаем с учётом равенства (12) при $p_T = r_T/l$

$$\sum_T l s_{Ti} = l \sum_{T: i \in T} s_{Ti} = l \sum_{T: i \in T} \frac{r_T q_i}{l \sum_{j \in T} q_j} = l q_i. \quad (17)$$

Возьмём набор z_{Ti} ($T \in \mathcal{T}, i \in M$) целых чисел, такой, что при $i \in T$ величины $|z_{Ti} - l s_{Ti}|$ ограничены некоторой общей константой и выполняются равенства

$$\sum_i z_{Ti} = r_T, \quad T \in \mathcal{T}, \quad (18)$$

а при $i \notin T$ величины $|z_{Ti}|$ равны 0. Это легко сделать в силу (16), определения (15) величин s_{Ti} и того, что мощность множеств T не превосходит константы $|A_0|$.

Введём величины

$$m_i = \sum_T z_{Ti}, \quad i \in M, \quad (19)$$

и положим $\mathbf{m} = (m_i, i \in M)$. Тогда

$$\sum_i m_i = \sum_{i,T} z_{Ti} = \sum_T \sum_i z_{Ti} = \sum_T r_T = l. \quad (20)$$

Кроме того, из (17), условий на числа s_{Ti} и ограниченности мощности множеств T константой $|\mathcal{T}|$ следует, что все величины $|m_i - lq_i|$ не превосходят некоторой общей константы.

Принимая во внимание (20), образуем частотный класс $\mathcal{K}_l(\mathbf{m})$ и используем его в качестве базисного множества в градиентной процедуре построения доопределений для слов класса $\mathcal{K}_l(\mathbf{r})$. Построим таблицу, строки которой соответствуют словам $v \in \mathcal{K}_l(\mathbf{r})$, столбцы — словам $w \in \mathcal{K}_l(\mathbf{m})$. На пересечении строки v и столбца w поместим 1 либо 0 в зависимости от того, доопределяет слово w слово v или нет. Числа b и d соответственно строк и столбцов этой таблицы задаются выражениями

$$b = \frac{l!}{\prod_T r_T!}, \quad d = \frac{l!}{\prod_i m_i!}.$$

При оценке снизу числа единиц в произвольной строке v таблицы будем учитывать только такие доопределения слова v , в которых для всех T и i символ a_T заменяется символом a_i в z_{Ti} позициях. Из (18) следует, что такие доопределения слова v существуют, а из (19) — что все они принадлежат классу $\mathcal{K}_l(\mathbf{m})$. Для каждого $v \in \mathcal{K}_l(\mathbf{r})$ число доопределений указанного вида равно

$$a = \prod_T \frac{r_T!}{\prod_i z_{Ti}!} = \frac{\prod_T r_T!}{\prod_{T,i} z_{Ti}!}.$$

Воспользовавшись утверждением 2 при указанных значениях b , d и a , а также тем, что $\frac{ba}{d} < b < l! < l^l$, оценим размер k градиентного покрытия таблицы:

$$k \leq \frac{l! \prod_{T,i} z_{Ti}!}{\prod_T r_T! \prod_i m_i!} l \ln l.$$

Комбинаторная энтропия $\log N_l(\mathbf{r})$ не превосходит $\log k$. Отсюда, применяя формулу Стирлинга к факториалам в оценке для k и учитывая равенства $\sum_{T,i} z_{Ti} = \sum_T r_T = \sum_i m_i = l$, приходим к оценке

$$\log N_l(\mathbf{r}) \leq l \log l + \sum_{T,i} z_{Ti} \log z_{Ti} - \sum_T r_T \log r_T - \sum_i m_i \log m_i + C_4 \log l. \quad (21)$$

Заменим в этом выражении величины z_{Ti} на ls_{Ti} , а m_i — на lq_i . Поскольку каждая из этих величин не превосходит l и заменяется величиной, отличающейся от неё не более

чем на константу, соответствующий ей член в (21) изменится не более чем на $C_5 \log l$. Число этих членов в (21) ограничено константой, а потому оценка (21) превращается в

$$\begin{aligned} \log N_l(\mathbf{r}) &\leq l \log l + \sum_{T,i} l s_{Ti} \log(l s_{Ti}) - \sum_T r_T \log r_T - \sum_i l q_i \log(l q_i) + C_3 \log l = \\ &= -l \left(\sum_{T,i} s_{Ti} \log s_{Ti} - \sum_T \frac{r_T}{l} \log \frac{r_T}{l} - \sum_i m_i \log q_i \right) + C_3 \log l. \end{aligned}$$

Принимая во внимание (16) и (17), получаем

$$\begin{aligned} \log N_l(\mathbf{r}) &\leq -l \left(\sum_{T,i} s_{Ti} \log s_{Ti} - \sum_{T,i} s_{Ti} \log \frac{r_T}{l} - \sum_{T,i} s_{Ti} \log q_i \right) + C_3 \log l = \\ &= l \sum_{T,i} s_{Ti} \log \frac{l s_{Ti}}{r_T q_i} + C_3 \log l. \end{aligned}$$

Замена s_{Ti} в аргументе логарифма значением из (15) приводит к оценке

$$\log N_l(\mathbf{r}) \leq l \sum_{T,i} s_{Ti} \log \frac{1}{\sum_{j \in T} q_j} + C_3 \log l = - \sum_T r_T \log \sum_{j \in T} q_j + C_3 \log l,$$

которая в силу (14) совпадает с (13). ■

4. Представительное множество системы частотных классов

Обозначим через $\mathfrak{K}_{\lambda, \mu}$ множество всех частотных классов $\mathcal{K}_l(\mathbf{r})$ (в рассматриваемом недоопределённом алфавите A), для которых длина l слов не превышает λ , а квазиэнтропия $h_l(\mathbf{r})$ не больше μ , где λ и μ — заданные натуральные параметры. Поскольку квазиэнтропия слов длины l не превосходит $l \log m$ (п. 1° леммы 2), можно считать

$$\mu \leq \lambda \log m, \quad (22)$$

а потому $\mu = O(\lambda)$.

Используемая в работе процедура кодирования недоопределённых источников требует нахождения доопределений для всех классов системы $\mathfrak{K}_{\lambda, \mu}$. При рассмотрении систем $\mathfrak{K}_{\lambda, \mu}$ в произвольных недоопределённых алфавитах возникают трудности, связанные с возможностью эффективной проверки условия $h_l(\mathbf{r}) \leq \mu$.

Для частично определённых алфавитов подобных трудностей не возникает, ибо для них квазиэнтропия выразима в явном виде (п. 7° леммы 2). В последующем ограничимся частично определёнными алфавитами, но развитый подход может быть распространён в некотором модифицированном виде и на другие алфавиты с явно выразимой энтропией (примеры таких алфавитов имеются в [5]). Дальше в данном пункте считаем, что $\mathfrak{K}_{\lambda, \mu}$ — система в алфавите $A = A_0 \cup \{*\} = \{a_0, a_1, \dots, a_{m-1}, *\}$. Она образована частотными классами $\mathcal{K}_l(\mathbf{r})$, $\mathbf{r} = (r_0, r_1, \dots, r_{m-1}, r_*)$, удовлетворяющими условиям $l \leq \lambda$, $h_l(\mathbf{r}) \leq \mu$.

Лемма 4. В случае частично определённых алфавитов условие $h(v) \leq \mu$ проверяемо со сложностью, ограниченной полиномом от длины $|v|$ слова v .

Доказательство. Поскольку параметр μ целочислен, условие $h(v) \leq \mu$ можно заменить на $\lceil h(v) \rceil \leq \mu$ ($\lceil x \rceil$ означает целое число, ближайшее к x сверху). В силу (8) оно приобретает вид $\lceil \log \alpha \rceil \leq \mu$, где

$$\alpha = \frac{(|v| - r_*(v))^{|v| - r_*(v)}}{\prod_{i=0}^{m-1} r_i(v)^{r_i(v)}}.$$

Нетрудно понять, что $\lceil \log \alpha \rceil$ совпадает с длиной двоичной записи числа $\lceil \alpha \rceil - 1$. Поэтому для проверки соотношения $h(v) \leq \mu$ достаточно найти по слову v двоичные записи параметров $r_i(v)$ и $r_*(v)$, вычислить по ним число $\lceil \alpha \rceil - 1$ в двоичной записи и сравнить длину этой записи с μ .

Для нахождения величины z^z , $z \in \{r_0(v), \dots, r_{m-1}(v), |v| - r_*(v)\}$, входящей в выражение для α , достаточно выполнить не более $2 \log z \leq 2 \log l$ умножений чисел, длина двоичных записей которых не превышает $l \log l$. Это требует не более $O(l^2 \log^3 l)$ (булевых) операций. Поскольку число значений z ограничено константой m' , на вычисление всех z^z также затрачивается не более $O(l^2 \log^3 l)$ операций. При известных значениях z^z число α может быть найдено применением $m - 1$ операций умножения и одной операции деления нацело (с избытком). На это расходуется $O(l^2 \log^2 l)$ операций. Общее число операций для вычисления α не превосходит $O(l^2 \log^3 l)$. Учёт других операций, используемых для проверки соотношения $h(v) \leq \mu$, не изменяет эту оценку. ■

Задача построения доопределений для всех классов системы $\mathfrak{K}_{\lambda, \mu}$ может быть сведена к аналогичной задаче для некоторой её подсистемы. Поставим задачу выделения в определённом смысле лучшей такой подсистемы. Отметим, что для доказательства основного результата работы (теоремы 2) достаточно более грубых рассмотрений, когда доопределения строятся для всех классов системы $\mathfrak{K}_{\lambda, \mu}$. Методы максимально возможного сокращения числа рассматриваемых классов могут быть полезны в приложениях.

Пусть π — перестановка множества $\{0, \dots, m - 1\}$. Для слова v (в алфавите $A = A_0 \cup \{*\}$) посредством πv будем обозначать слово, полученное из v заменой символов $a_i \in A_0$ на $a_{\pi(i)}$ и сохранением символов $*$. Для множества слов V положим $\pi V = \{\pi v, v \in V\}$. Если V — некоторое множество слов с длиной не меньше l , то через $V|_l$ будем обозначать множество начал длины l всех слов из V .

Пусть $\mathcal{K}_{l_1}(\mathbf{r}_1)$ и $\mathcal{K}_{l_2}(\mathbf{r}_2)$ — частотные классы слов (в частично определённом алфавите A). Будем говорить, что класс $\mathcal{K}_{l_1}(\mathbf{r}_1)$ *представительнее* класса $\mathcal{K}_{l_2}(\mathbf{r}_2)$, если $l_1 \geq l_2$ и найдётся перестановка π множества $\{0, \dots, m - 1\}$, такая, что каково бы ни было доопределение $\mathcal{D}_{l_1}(\mathbf{r}_1)$ класса $\mathcal{K}_{l_1}(\mathbf{r}_1)$, множество $(\pi \mathcal{D}_{l_1}(\mathbf{r}_1))|_{l_2}$ образует доопределение класса $\mathcal{K}_{l_2}(\mathbf{r}_2)$. Заметим, что если в определении представительности вместо начал слов брать другие подслова, расположенные в l_2 различных фиксированных разрядах, это не изменит отношения представительности, поскольку частотные классы замкнуты относительно перестановок символов в словах.

Подсистему \mathfrak{M} некоторой системы \mathfrak{K} частотных классов назовём *представительным множеством системы \mathfrak{K}* , если для любого класса $\mathcal{K}_l(\mathbf{r}) \in \mathfrak{K}$ в \mathfrak{M} имеется более представительный класс. Таким образом, доопределение любого класса системы \mathfrak{K} может быть образовано из доопределения подходящего класса системы \mathfrak{M} переименованием символов алфавита и отбрасыванием некоторого количества заключительных символов слов. Представительное множество \mathfrak{M} называется *минимальным*, если не существует представительного множества системы \mathfrak{K} , содержащего меньшее число частотных классов. Целью данного пункта является явное описание некоторого минимального представительного множества для системы $\mathfrak{K}_{\lambda, \mu}$.

Наряду с наборами $\mathbf{r} = (r_0, \dots, r_{m-1}, r_*)$ будем рассматривать соответствующие укороченные наборы $\hat{\mathbf{r}} = (r_0, \dots, r_{m-1})$. Будем говорить, что укороченный набор $\hat{\mathbf{r}}_1 = (r_{1,0}, \dots, r_{1,m-1})$ *мажорирует* $\hat{\mathbf{r}}_2 = (r_{2,0}, \dots, r_{2,m-1})$, и записывать $\hat{\mathbf{r}}_1 \geq \hat{\mathbf{r}}_2$, если $r_{1,0} \geq r_{2,0}, \dots, r_{1,m-1} \geq r_{2,m-1}$. Величину $W(\hat{\mathbf{r}}) = r_0 + \dots + r_{m-1}$ назовём *весом* укороченного набора $\hat{\mathbf{r}}$.

ченного набора $\hat{\mathbf{r}}$. С набором $\hat{\mathbf{r}}$ свяжем его квазиэнтропию

$$h(\hat{\mathbf{r}}) = W(\hat{\mathbf{r}}) \log W(\hat{\mathbf{r}}) - \sum_{i=0}^{m-1} r_i \log r_i.$$

В силу (11) квазиэнтропия $h_l(\mathbf{r})$ класса $\mathcal{K}_l(\mathbf{r})$ при любом $l \geq W(\hat{\mathbf{r}})$ совпадает с квазиэнтропией $h(\hat{\mathbf{r}})$ укороченного набора.

Результатом применения перестановки π к укороченному набору $\hat{\mathbf{r}}$ считается набор $\pi\hat{\mathbf{r}} = (r_{\pi(0)}, \dots, r_{\pi(m-1)})$. Набор $\hat{\mathbf{r}}$, для которого $r_0 \geq \dots \geq r_{m-1}$, будем называть *упорядоченным*. Перестановка π , переводящая заданный укороченный набор $\hat{\mathbf{r}}$ в упорядоченный набор $\pi\hat{\mathbf{r}}$, не единственна, ибо разным i могут соответствовать одинаковые r_i . Для её конкретизации дополнительно потребуем, чтобы при $r_i = r_j$ и $i < j$ было выполнено $\pi(i) < \pi(j)$. Эту перестановку обозначим $\pi_{\hat{\mathbf{r}}}$ и назовём *упорядочивающей* (для $\hat{\mathbf{r}}$). Для упорядоченных наборов $\hat{\mathbf{r}}_1$ и $\hat{\mathbf{r}}_2$ будем говорить, что $\hat{\mathbf{r}}_2$ *непосредственно следует* за $\hat{\mathbf{r}}_1$ ($\hat{\mathbf{r}}_1$ *непосредственно предшествует* $\hat{\mathbf{r}}_2$), если $\hat{\mathbf{r}}_2 \geq \hat{\mathbf{r}}_1$ и $W(\hat{\mathbf{r}}_2) = W(\hat{\mathbf{r}}_1) + 1$. Так, например, всеми непосредственно следующими за набором $(3,1,1,0,0)$ являются наборы $(4,1,1,0,0)$, $(3,2,1,0,0)$ и $(3,1,1,1,0)$.

Эффективный метод проверки наличия отношения представительности для заданных частотных классов даётся следующим утверждением.

Утверждение 4. Класс $\mathcal{K}_{l_1}(\mathbf{r}_1)$ представительнее класса $\mathcal{K}_{l_2}(\mathbf{r}_2)$ тогда и только тогда, когда $l_1 \geq l_2$ и набор $\pi_{\hat{\mathbf{r}}_1}\hat{\mathbf{r}}_1$, полученный упорядочением набора $\hat{\mathbf{r}}_1$, мажорирует набор $\pi_{\hat{\mathbf{r}}_2}\hat{\mathbf{r}}_2$, полученный упорядочением $\hat{\mathbf{r}}_2$.

Доказательство. Будем считать $l_1 \geq l_2$.

1. Пусть имеет место мажорирование $\pi_{\hat{\mathbf{r}}_1}\hat{\mathbf{r}}_1 \geq \pi_{\hat{\mathbf{r}}_2}\hat{\mathbf{r}}_2$. Применив к обеим частям перестановку $\pi_{\hat{\mathbf{r}}_2}^{-1}$ и положив $\pi = \pi_{\hat{\mathbf{r}}_2}^{-1}\pi_{\hat{\mathbf{r}}_1}$, получаем соотношение $\pi\hat{\mathbf{r}}_1 \geq \hat{\mathbf{r}}_2$. Обозначим через \mathbf{r}_2^+ результат дописывания $l_1 - l_2$ нулей к набору \mathbf{r}_2 . Из предшествующего соотношения и равенства $\hat{\mathbf{r}}_2^+ = \hat{\mathbf{r}}_2$ следует $\pi\hat{\mathbf{r}}_1 \geq \hat{\mathbf{r}}_2^+$.

Рассмотрим произвольное доопределение $\mathcal{D}_{l_1}(\mathbf{r}_1)$ класса $\mathcal{K}_{l_1}(\mathbf{r}_1)$. Ясно, что $\pi\mathcal{D}_{l_1}(\mathbf{r}_1)$ является некоторым доопределением класса $\mathcal{K}_{l_1}(\pi\mathbf{r}_1)$. Отсюда и из $\pi\hat{\mathbf{r}}_1 \geq \hat{\mathbf{r}}_2^+$ нетрудно заключить, что множество $\pi\mathcal{D}_{l_1}(\mathbf{r}_1)$ доопределяет класс $\mathcal{K}_{l_1}(\mathbf{r}_2^+)$, а потому $(\pi\mathcal{D}_{l_1}(\mathbf{r}_1))|_{l_2}$ образует доопределение класса $\mathcal{K}_{l_2}(\mathbf{r}_2)$. Это означает, что класс $\mathcal{K}_{l_1}(\mathbf{r}_1)$ представительнее $\mathcal{K}_{l_2}(\mathbf{r}_2)$.

2. Обратное, пусть $\mathcal{K}_{l_1}(\mathbf{r}_1)$ представительнее $\mathcal{K}_{l_2}(\mathbf{r}_2)$. Перестановка π , присутствующая в определении соотношения этих классов по представительности, удовлетворяет условию $\pi\hat{\mathbf{r}}_1 \geq \hat{\mathbf{r}}_2$. Действительно, если это не так, то при некотором i имеет место $r_{1,\pi^{-1}(i)} < r_{2,i}$. Тогда, используя для слов класса $\mathcal{K}_{l_1}(\mathbf{r}_1)$ в качестве доопределений символов * символы, отличные от $a_{\pi^{-1}(i)}$, можно получить доопределение $\mathcal{D}_{l_1}(\mathbf{r}_1)$, все слова которого содержат менее $r_{2,i}$ символов $a_{\pi^{-1}(i)}$, а потому число символов a_i в любом слове множества $(\pi\mathcal{D}_{l_1}(\mathbf{r}_1))|_{l_2}$ меньше $r_{2,i}$. Это противоречит тому, что данное множество доопределяет класс $\mathcal{K}_{l_2}(\mathbf{r}_2)$, слова которого используют $r_{2,i}$ символов a_i .

Покажем теперь, что если для некоторых числовых наборов $\mathbf{s}_1 = (s_{1,0}, \dots, s_{1,m-1})$ и $\mathbf{s}_2 = (s_{2,0}, \dots, s_{2,m-1})$ выполнено условие мажорирования $\mathbf{s}_1 \geq \mathbf{s}_2$, то аналогичному условию мажорирования удовлетворяют и соответствующие им упорядоченные наборы. Для доказательства можно ограничиться случаем, когда набор \mathbf{s}_1 с самого начала является упорядоченным. Возьмём максимальный по величине разряд набора \mathbf{s}_2 (если таких разрядов несколько, то первый из них); пусть это будет $s_{2,i}$. Поменяем в \mathbf{s}_2 местами разряды $s_{2,1}$ и $s_{2,i}$. Так как $s_{2,i} \leq s_{1,i} \leq s_{1,1}$ и $s_{2,1} \leq s_{2,i} \leq s_{1,i}$, то полученный в результате набор \mathbf{s}'_2 будет удовлетворять условию мажорирования $\mathbf{s}_1 \geq \mathbf{s}'_2$. Аналогич-

ные соображения показывают, что если среди разрядов набора \mathbf{s}'_2 , отличных от первого, найти максимальный по величине и поменять его местами со вторым разрядом, получим набор \mathbf{s}''_2 , для которого $\mathbf{s}_1 \geq \mathbf{s}''_2$. Эти рассуждения могут быть продолжены вплоть до перевода набора \mathbf{s}_2 в упорядоченный.

Используя в качестве \mathbf{s}_1 и \mathbf{s}_2 наборы $\pi \hat{\mathbf{r}}_1$ и $\hat{\mathbf{r}}_2$, удовлетворяющие условию $\pi \hat{\mathbf{r}}_1 \geq \hat{\mathbf{r}}_2$, и учитывая, что наборам $\pi \hat{\mathbf{r}}_1$ и $\hat{\mathbf{r}}_1$ соответствуют одинаковые упорядоченные наборы, приходим к $\pi_{\hat{\mathbf{r}}_1} \hat{\mathbf{r}}_1 \geq \pi_{\hat{\mathbf{r}}_2} \hat{\mathbf{r}}_2$. ■

Набор $\hat{\mathbf{r}}$ назовём (λ, μ) -допустимым, если он упорядочен и удовлетворяет условиям $W(\hat{\mathbf{r}}) \leq \lambda$, $h(\hat{\mathbf{r}}) \leq \mu$. Упорядоченный набор $\hat{\mathbf{r}}$, мажорируемый некоторым (λ, μ) -допустимым набором $\hat{\mathbf{r}}_1$, также (λ, μ) -допустим. Действительно, произвольное слово v класса $\mathcal{K}_{W(\hat{\mathbf{r}})}(\hat{\mathbf{r}})$ может быть продолжено до некоторого слова v_1 класса $\mathcal{K}_{W(\hat{\mathbf{r}}_1)}(\hat{\mathbf{r}}_1)$. В силу п. 5° леммы 2 выполнено $h(v) \leq h(v_1)$, а потому из соотношений $h(v) = h_{W(\hat{\mathbf{r}})}(\hat{\mathbf{r}}) = h(\hat{\mathbf{r}})$, $h(v_1) = h_{W(\hat{\mathbf{r}}_1)}(\hat{\mathbf{r}}_1) = h(\hat{\mathbf{r}}_1)$ и из (λ, μ) -допустимости набора $\hat{\mathbf{r}}_1$ следует $h(\hat{\mathbf{r}}) \leq h(\hat{\mathbf{r}}_1) \leq \mu$. Неравенство $W(\hat{\mathbf{r}}) \leq W(\hat{\mathbf{r}}_1) \leq \lambda$ очевидно.

Отношение мажорирования $\hat{\mathbf{r}}_1 \geq \hat{\mathbf{r}}_2$ является частичным порядком. Обозначим через $R_{\lambda, \mu}$ множество всех (λ, μ) -допустимых наборов, максимальных по этому порядку. В нем отсутствуют мажорирования одних наборов другими и всякий (λ, μ) -допустимый набор мажорируется некоторым набором из этого множества.

Опишем прямой способ построения множества $R_{\lambda, \mu}$. Сначала индуктивно построим последовательность множеств $R'_{\lambda, \mu}(i)$, $i = 0, \dots, \lambda$. Для этого положим $R'_{\lambda, \mu}(0) = \{(0, 0, \dots, 0)\}$, а множество $R'_{\lambda, \mu}(i+1)$ образуем, взяв все наборы, непосредственно следующие за наборами из $R'_{\lambda, \mu}(i)$, и устранив те из них, квазиэнтропия которых превышает μ . Обозначим через $R_{\lambda, \mu}(i)$ ($i = 0, \dots, \lambda$) результат удаления из $R'_{\lambda, \mu}(i)$ наборов, мажорируемых какими-либо наборами множества $R'_{\lambda, \mu}(i+1)$ (считаем $R'_{\lambda, \mu}(\lambda+1) = \emptyset$).

Утверждение 5. Множество $R_{\lambda, \mu}(i)$ ($i = 0, 1, \dots, \lambda$) образовано всеми наборами множества $R_{\lambda, \mu}$, имеющими вес i , а потому $R_{\lambda, \mu} = \bigcup_{i=0}^{\lambda} R_{\lambda, \mu}(i)$. Сложность построения множества $R_{\lambda, \mu}$ полиномиальна по λ .

Доказательство. Сначала индукцией по i убедимся, что множество $R'_{\lambda, \mu}(i)$ состоит из всех (λ, μ) -допустимых наборов веса i . Для $R'_{\lambda, \mu}(0)$ это очевидно. Считая факт справедливым для $R'_{\lambda, \mu}(i)$, $i < \lambda$, рассмотрим множество $R'_{\lambda, \mu}(i+1)$. Оно непусто, поскольку включает набор $(1, \dots, 1, 0, \dots, 0)$ веса $i+1$, имеющий нулевую квазиэнтропию. Наборы этого множества (λ, μ) -допустимы, ибо, непосредственно следуя за наборами из $R'_{\lambda, \mu}(i)$, они имеют вес $i+1 \leq \lambda$, а их квазиэнтропия не превышает μ по построению. С другой стороны, если некоторый набор $\hat{\mathbf{r}}$ веса $i+1$ (λ, μ) -допустим, то любой непосредственно предшествующий ему набор также (λ, μ) -допустим, а его вес равен i . По предположению индукции он содержится в $R'_{\lambda, \mu}(i)$, а потому непосредственно следующий за ним набор $\hat{\mathbf{r}}$ включается в $R'_{\lambda, \mu}(i+1)$ по построению. Индукция завершена.

Покажем теперь, что $R_{\lambda, \mu}(i)$ — множество всех максимальных по мажорированию (λ, μ) -допустимых наборов веса i . Произвольный (λ, μ) -допустимый набор $\hat{\mathbf{r}} = (r_0, \dots, r_{m-1})$ веса i содержится в $R'_{\lambda, \mu}(i)$, и если $\hat{\mathbf{r}}$ не мажорируется, то, в частности, он не мажорируется и наборами из $R'_{\lambda, \mu}(i+1)$, а потому будет включён в $R_{\lambda, \mu}(i)$. Предположим теперь, что $\hat{\mathbf{r}}$ мажорируется некоторым (λ, μ) -допустимым набором $\hat{\mathbf{r}}' = (r'_0, \dots, r'_{m-1})$, и пусть эти наборы впервые различаются в разряде j , т. е. $r_0 = r'_0, \dots, r_{j-1} = r'_{j-1}$ и $r_j < r'_j$. образуем набор $\hat{\mathbf{r}}'' = (r_0, \dots, r_{j-1}, r_j + 1, r_{j+1}, \dots, r_{m-1})$. Он упорядочен и мажорируется (λ, μ) -допустимым набором $\hat{\mathbf{r}}'$, а потому (λ, μ) -допу-

стим. Этот набор имеет вес $i + 1$ и, следовательно, содержится в $R'_{\lambda,\mu}(i + 1)$. Набор $\hat{\mathbf{r}}''$ мажорирует $\hat{\mathbf{r}}$, и при построении множества $R_{\lambda,\mu}(i)$ набор $\hat{\mathbf{r}}$ будет исключен из $R'_{\lambda,\mu}(i)$.

Тот факт, что сложность построения множества $R_{\lambda,\mu}$ указанным способом полиномиальна по λ , легко извлекается из леммы 4 и того, что число укороченных наборов $\hat{\mathbf{r}}$ веса $W(\hat{\mathbf{r}}) \leq \lambda$ ограничено полиномом от λ . ■

С каждым набором $\hat{\mathbf{r}} \in R_{\lambda,\mu}$ свяжем частотный класс $\mathcal{K}_\lambda(\mathbf{r})$, где набор \mathbf{r} получен из $\hat{\mathbf{r}}$ дописыванием компоненты $r_* = \lambda - W(\hat{\mathbf{r}})$. Совокупность этих частотных классов обозначим через $\mathfrak{M}_{\lambda,\mu}$.

Утверждение 6. Система $\mathfrak{M}_{\lambda,\mu}$ образует минимальное представительное множество системы $\mathfrak{K}_{\lambda,\mu}$. Её явное задание может быть найдено с полиномиальной по λ сложностью.

Доказательство. Убедимся сначала, что $\mathfrak{M}_{\lambda,\mu}$ является представительным множеством системы $\mathfrak{K}_{\lambda,\mu}$. Рассмотрим произвольный частотный класс $\mathcal{K}_l(\mathbf{r}) \in \mathfrak{K}_{\lambda,\mu}$. Пусть $\hat{\mathbf{r}}$ — укороченный набор для \mathbf{r} , $\pi_{\hat{\mathbf{r}}}$ — упорядочивающая его перестановка. Легко видеть, что набор $\pi_{\hat{\mathbf{r}}}\hat{\mathbf{r}}$, полученный упорядочиванием $\hat{\mathbf{r}}$, (λ, μ) -допустим, а потому в $R_{\lambda,\mu}$ имеется мажорирующий его набор $\hat{\mathbf{r}}'$. Возьмём соответствующий ему класс $\mathcal{K}_\lambda(\mathbf{r}')$ системы $\mathfrak{M}_{\lambda,\mu}$. Нетрудно понять, что для любого доопределения $\mathcal{D}_\lambda(\mathbf{r}')$ этого класса множество $(\pi_{\hat{\mathbf{r}}}^{-1}\mathcal{D}_\lambda(\mathbf{r}'))|_l$ доопределяет класс $\mathcal{K}_l(\mathbf{r})$. Это означает, что класс $\mathcal{K}_\lambda(\mathbf{r}')$ представительнее $\mathcal{K}_l(\mathbf{r})$, а потому $\mathfrak{M}_{\lambda,\mu}$ — представительное множество системы $\mathfrak{K}_{\lambda,\mu}$.

Покажем теперь, что представительное множество $\mathfrak{M}_{\lambda,\mu}$ минимально. Предположим, что это не так и существует представительное множество \mathfrak{M}' меньшей мощности. Тогда в \mathfrak{M}' найдётся класс $\mathcal{K}_\lambda(\mathbf{r}')$, который представительнее по крайней мере двух различных классов $\mathcal{K}_\lambda(\mathbf{r}_1)$ и $\mathcal{K}_\lambda(\mathbf{r}_2)$ системы $\mathfrak{M}_{\lambda,\mu}$. Наборы $\hat{\mathbf{r}}_1$ и $\hat{\mathbf{r}}_2$ упорядочены по определению и, согласно утверждению 4, мажорируются результатом упорядочивания набора $\hat{\mathbf{r}}'$. Будучи максимальными по отношению мажорирования, они совпадают с ним и, следовательно, одинаковы. Это противоречит тому, что классы $\mathcal{K}_\lambda(\mathbf{r}_1)$ и $\mathcal{K}_\lambda(\mathbf{r}_2)$ выбраны различными.

Система $\mathfrak{M}_{\lambda,\mu}$ явно задаётся множеством пар (\mathbf{r}, λ) , где \mathbf{r} — результат приписывания к набору $\hat{\mathbf{r}} \in R_{\lambda,\mu}$ компоненты $\lambda - W(\hat{\mathbf{r}})$. Согласно утверждению 5, сложность задания множества $R_{\lambda,\mu}$ полиномиальна по λ , а потому и сложность задания системы $\mathfrak{M}_{\lambda,\mu}$ полиномиальна по λ . ■

5. Справочная часть кода

Кодирование использует параметры λ и μ , которые будут назначены позже. Для каждого $v \in A^n$ кодовое слово $K(v)$ образовано подсловами K_0 и $K_1(v)$, $K(v) = K_0 K_1(v)$, где K_0 , одинаковое для всех v , называется *справочной частью* кодового слова, $K_1(v)$ — его *основной частью*. Справочная часть K_0 содержит сведения о параметре m задачи, параметрах λ и μ алгоритма кодирования, а также о доопределениях слов из частотных классов $\mathcal{K}_\lambda(\mathbf{r})$, входящих в минимальное представительное множество $\mathfrak{M}_{\lambda,\mu}$.

Будем использовать следующее представление натуральных чисел s в виде двоичных слов \hat{s} . Если $\sigma_1 \dots \sigma_k$, $k = \lceil \log(s - 1) \rceil$, — минимальная двоичная запись числа s , то $\hat{s} = \sigma_1 \sigma_1 \dots \sigma_k \sigma_k 01$. Такое представление чисел удобно применять для задания совокупностей чисел и двоичных слов. Так, например, если s — число, а u_1 и u_2 — слова, то слово $\hat{s}u_1$ однозначно задаёт пару (s, u_1) , а слово $\widehat{|u_1|}u_1u_2$, где $\widehat{|u_1|}$ — представление длины $|u_1|$ слова u_1 , однозначно задаёт пару слов (u_1, u_2) . Очевидно, что

$$|\hat{s}| \leq 2 \log s + 4. \quad (23)$$

Для каждого класса $\mathcal{K}_\lambda(\mathbf{r})$ построим доопределение $\mathcal{D}_\lambda(\mathbf{r})$ с использованием градиентной процедуры из леммы 3. Обозначим через $\mathcal{D}_{\lambda,\mu}$ объединение доопределений $\mathcal{D}_\lambda(\mathbf{r})$ для всех классов $\mathcal{K}_\lambda(\mathbf{r})$ системы $\mathfrak{M}_{\lambda,\mu}$. Множество $\mathcal{D}_{\lambda,\mu}$ обладает тем свойством, что для любого слова v с $|v| \leq \lambda$ и $h(v) \leq \mu$ в нём найдётся слово, начало которого доопределяет слово v' , полученное из v упорядочивающей перестановкой. Множество с таким свойством будем называть (λ, μ) -представительным. Множество $\mathcal{D}_{\lambda,\mu}$ упорядочим лексикографически.

Слова $w \in \mathcal{D}_{\lambda,\mu}$ в алфавите $A_0 = \{a_0, \dots, a_{m-1}\}$ имеют длину λ . Сопоставим каждому w двоичное слово \tilde{w} длины $m\lambda$, заменив в w символы a_i двоичными словами $0\dots 010\dots 0$ длины m , содержащими 1 в разряде i , $0 \leq i \leq m-1$. Слова \tilde{w} также будут расположены в лексикографическом порядке в соответствии с упорядочением слов w . образуем слово \tilde{w}_Σ путём приписывания в этом порядке слов \tilde{w} друг к другу. В качестве справочной части кодовых слов $K(v)$, $v \in A^n$ будем использовать двоичное слово

$$K_0 = \widehat{m\lambda\mu}|\widehat{\tilde{w}_\Sigma}|\tilde{w}_\Sigma. \quad (24)$$

Подслово $|\widehat{\tilde{w}_\Sigma}|$ позволяет при декодировании кодового слова $K(v) = K_0 K_1(v)$ отделить в нём справочную часть от основной.

Лемма 5.

- 1) Длина справочной части кода удовлетворяет оценке

$$|K_0| \leq \lambda^{C_6} 2^\mu. \quad (25)$$

- 2) Сложность построения справочной части ограничена величиной $(C_7)^\lambda$.

Здесь C_6 и C_7 — подходящие константы.

Доказательство.

- 1) По лемме 3 с учётом неравенства $h_\lambda(\mathbf{r}) \leq \mu$ получаем оценку мощности множества $\mathcal{D}_\lambda(\mathbf{r})$, доопределяющего класс $\mathcal{K}_\lambda(\mathbf{r})$:

$$|\mathcal{D}_\lambda(\mathbf{r})| \leq \lambda^{C_3} 2^\mu.$$

Для мощности объединения $\mathcal{D}_{\lambda,\mu}$ этих множеств справедливы оценки

$$|\mathcal{D}_{\lambda,\mu}| \leq \lambda^{C_3} 2^\mu |\mathfrak{M}_{\lambda,\mu}| \leq \lambda^{C_3} 2^\mu \lambda^{m+1}, \quad (26)$$

и потому длина слова \tilde{w}_Σ удовлетворяет оценке

$$|\tilde{w}_\Sigma| \leq m\lambda^{C_3+m+1} 2^\mu.$$

Отсюда, учитывая (22), (23) и (24), приходим к оценке (25).

- 2) Оценим сложность градиентной процедуры доопределения класса $\mathcal{K}_\lambda(\mathbf{r}) \in \mathfrak{M}_{\lambda,\mu}$. Необходимый для построения градиентной таблицы набор Q , на котором достигается квазиэнтропия $h_\lambda(\mathbf{r})$, вычисляется с полиномиальной сложностью по формуле $Q = \mathbf{r}/(\lambda - r_*)$ (лемма 2, п. 7°). Числа строк и столбцов градиентной таблицы не более чем экспоненциальны (относительно λ), поэтому по замечанию 3 процедура реализуется с экспоненциальной оценкой сложности. Поскольку множество $\mathfrak{M}_{\lambda,\mu}$ содержит полиномиальное число классов $\mathcal{K}_\lambda(\mathbf{r})$, сложность построения доопределений всех слов из классов множества $\mathfrak{M}_{\lambda,\mu}$ остаётся экспоненциальной. Нетрудно понять, что учёт других операций, используемых для построения справочной части кода, не изменяет характер общей оценки сложности. ■

6. Основная часть кода

Перед построением основной части кодового слова производится обработка справочной части $K_0 = \widehat{m}\widehat{\lambda}\widehat{\mu}|\widehat{w}_\Sigma|\widehat{w}_\Sigma$. По ней находятся параметры m , λ , μ и слово \tilde{w}_Σ . Затем слово \tilde{w}_Σ разбивается на подслова \tilde{w} длины $m\lambda$ и по ним восстанавливаются соответствующие слова w длины λ в алфавите A_0 . Согласно (26), число слов w не превосходит $\lambda^{C_8}2^\mu$. Слова w нумеруются в лексикографическом порядке двоичными числами (т. е. числами в двоичной записи) $\tilde{\delta}(w)$ длины $d = \lceil \log(\lambda^{C_8}2^\mu) \rceil$. Кроме того, все перестановки π множества $\{0, 1, \dots, m-1\}$ нумеруются двоичными числами $\tilde{\varepsilon}(\pi)$ длины $z = \lceil \log m! \rceil$ (напомним, что m — константа рассматриваемой задачи).

Опишем метод построения основной части $K_1(v)$ кода для слова $v \in A^n$. Разобьём слово v на куски v_i , $i = 1, 2, \dots$, последовательно отрезая от него подслова v_i максимально возможной длины $|v_i|$, удовлетворяющие условиям $h(v_i) \leq \mu$ и $|v_i| \leq \lambda$. Если число этих кусков равно t , то $v = v_1v_2 \dots v_t$. Такое разбиение однозначно в силу монотонности квазиэнтропии (п. 5° леммы 2). Положим $b = \lceil \log \lambda \rceil$ и b -разрядную двоичную запись длины куска v_i обозначим через $\tilde{\lambda}_i$.

Каждому куску v_i сопоставим его код $K(v_i)$ следующим образом. Рассмотрим укороченный набор $\mathbf{r}(v_i) = (r_0(v_i), \dots, r_{m-1}(v_i))$ для слова v_i , построим по нему упорядочивающую перестановку $\pi_{\mathbf{r}(v_i)}$ и образуем слово $\pi_{\mathbf{r}(v_i)}v_i = v'_i$. Слова w , возникшие в результате обработки справочной части кода, образуют (λ, μ) -представительное множество, поэтому в нём имеются слова, начала которых доопределяют v'_i . Возьмём любое из них, обозначим его через $w^{(i)}$ и положим

$$K(v_i) = \tilde{\delta}(w^{(i)})\tilde{\lambda}_i\tilde{\varepsilon}(\pi_{\mathbf{r}(v_i)}^{-1}). \quad (27)$$

В качестве основной части кода для слова $v \in A^n$ возьмём

$$K_1(v) = \widehat{d}\widehat{b}\widehat{z}\widehat{t}K(v_1)K(v_2) \dots K(v_t). \quad (28)$$

Из этих выражений с учётом (23) следует, что длина основной части кода может быть оценена величиной

$$|K_1(v)| \leq 2 \log d + 2 \log b + 2 \log z + 2 \log t + t(d + b + z) + 16. \quad (29)$$

Лемма 6.

1) Длина основной части $K_1(v)$ кода для слова $v \in A^n$ удовлетворяет оценке

$$|K_1(v)| \leq h(v) + n \left(\frac{C_9 \log \lambda}{\mu - \log \lambda - 2} + \frac{\mu + C_9 \log \lambda}{\lambda} \right) + 2 \log n + C_{10} \lambda, \quad (30)$$

где $h(v)$ — квазиэнтропия слова v .

2) Сложность построения основной части кода для слова $v \in A^n$ при заданной проверочной части ограничена величиной $n\lambda^{C_{11}}2^\mu$.

Доказательство.

1) Будем говорить, что кусок v_i слова v имеет тип 1, если $|v_i| < \lambda$, и тип 2, если $|v_i| = \lambda$. Количества кусков типа 1 и 2 в слове v обозначим соответственно t_1 и t_2 . Очевидно, что $t_2 \leq n/\lambda$.

Оценим t_1 . Если кусок v_i имеет тип 1, то либо он является заключительным в слове v , либо квазиэнтропия $h(v_i a)$ слова, полученного приписыванием к v_i символа a ,

расположенного в слове v после v_i , превышает μ . Во втором случае в силу п. 6° леммы 2 имеем

$$h(v_i) \geq h(v_i a) - \log |v_i| - \log e \geq \mu - \log \lambda - 2.$$

Используя это неравенство и тот факт, что число незаключительных кусков типа 1 не превышает $t_1 - 1$, выводим на основе п. 4° леммы 2

$$h(v) = h(v_1 \dots v_t) \geq \sum_{i=1}^t h(v_i) \geq (t_1 - 1)(\mu - \log \lambda - 2).$$

Отсюда $t_1 \leq \frac{h(v)}{\mu - \log \lambda - 2} + 1$, а потому $t = t_1 + t_2 \leq \frac{h(v)}{\mu - \log \lambda - 2} + \frac{n}{\lambda} + 1$.

Перепишем (29) в виде

$$|K_1(v)| \leq [(t_1 - 1)(d + b + z)] + [t_2(d + b + z)] + [(d + b + z) + 2(\log d + \log b + \log z) + 2 \log n + 16]$$

и найдём оценку этого выражения как сумму оценок трёх его фрагментов, выделенных квадратными скобками. Имеем

$$\begin{aligned} (t_1 - 1)(d + b + z) &\leq \frac{h(v)}{\mu - \log \lambda - 2} ((C_8 + 1) \log \lambda + \mu + m \log m + 3) = \\ &= h(v) + \frac{h(v)}{\mu - \log \lambda - 2} ((C_8 + 2) \log \lambda + m \log m + 5). \end{aligned}$$

Оценив здесь второе вхождение $h(v)$ максимально возможным значением $n \log m$ квазиэнтропии (п. 1° леммы 2), получаем

$$\begin{aligned} (t_1 - 1)(d + b + z) &\leq h(v) + \frac{n \log m ((C_8 + 2) \log \lambda + m \log m + 5)}{\mu - \log \lambda - 2} \leq \\ &\leq h(v) + \frac{C_9 n \log \lambda}{\mu - \log \lambda - 2}. \end{aligned} \quad (31)$$

Оценка второго фрагмента даёт

$$t_2(d + b + z) \leq \frac{n((C_8 + 1) \log \lambda + \mu + m \log m + 3)}{\lambda} \leq n \frac{\mu + C_9 \log \lambda}{\lambda}. \quad (32)$$

Оценим заключительный фрагмент, учитывая неравенства $t \leq n$ и $\mu \leq \lambda \log m$:

$$\begin{aligned} (d + b + z) + 2(\log d + \log b + \log z) + 2 \log t + 16 &\leq \\ &\leq (d + b + z) + 6 \log(d + b + z) + 2 \log t + 16 \leq \\ &\leq \mu + C_9 \log \lambda + 6 \log(\mu + C_9 \log \lambda) + 2 \log n + 16 \leq 2 \log n + C_{10} \lambda. \end{aligned} \quad (33)$$

Просуммировав (31), (32) и (33), получаем (30).

2) Алгоритм построения основной части $K_1(v)$ кода при заданной справочной части K_0 включает три основных этапа.

На первом этапе производится обработка справочной части кода, как описано в начале раздела. С учётом оценки (25) нетрудно заключить, что сложность этого этапа не превышает $\lambda^{C_{12}} 2^\mu$.

Второй этап посвящён разбиению слова v на куски v_i . Сложность этого этапа ограничена величиной $n \lambda^{C_{13}}$, где n — длина слова v .

Третий этап включает построение кодов $K(v_i)$ всех кусков в соответствии с (27) и образование из них основной части кода для слова v . Наибольшая сложность при построении $K(v_i)$ приходится на поиск слова $w^{(i)}$, начало которого доопределяет слово v'_i , полученное из v_i упорядочивающей перестановкой. Слово $w^{(i)}$ может быть найдено путём последовательного просмотра слов w , возникших в результате обработки справочной части, и проверки, доопределяют ли их начала слово v'_i . С учётом того, что число слов w не превосходит $\lambda^{C_8} 2^\mu$, для этого достаточно выполнить не более $\lambda^{C_{14}} 2^\mu$ операций. Можно считать, что эта оценка покрывает также полиномиальное по λ число других операций, требуемых для построения кода $K(v_i)$ в соответствии с (27). Поскольку число t кусков v_i не превышает длину n слова v , на построение кодов всех кусков затрачивается не более $n\lambda^{C_{14}} 2^\mu$ операций. Будем считать, что константа C_{14} выбрана так, что эта оценка учитывает также операции, используемые для построения $K_1(v)$ из кодов $K(v_i)$ в соответствии с (28) и, таким образом, оценивает сложность третьего этапа.

Величина сложности из п. 2 формулировки леммы 6 оценивает сверху при подходящем выборе C_{11} сумму сложностей всех трёх этапов. ■

7. Оценка характеристик кодирования

Покажем, что предложенное кодирование K удовлетворяет исходным требованиям к кодированию частично определённых слов.

Лемма 7. Кодирование K разделимо и позволяет для всякого слова $v \in A^n$ восстановить по его коду $K(v)$ некоторое доопределение слова v .

Доказательство.

1. Убедимся, что код обладает достаточным для делимости свойством префикса. Пусть $K(v)$ и $K(v')$ — кодовые слова и $K(v)$ является префиксом (началом) слова $K(v')$. Это соотношение переносится на их основные части $K_1(v)$ и $K_1(v')$, поскольку справочные части совпадают. Тогда из (28) следует, что $t' = t$ и слово $K(v_1) \dots K(v_t)$ является префиксом слова $K(v'_1) \dots K(v'_t)$. Так как эти слова имеют одинаковую длину $t(d + b + z)$, они обязаны совпасть, что влечёт совпадение кодовых слов $K(v)$ и $K(v')$.

2. Метод декодирования — восстановления по $K(v)$ некоторого доопределения слова v — включает два этапа.

На первом этапе сначала по начальной части $\widehat{m\lambda\mu|\widehat{w}_\Sigma}$ кодового слова $K(v)$ находятся параметры m , λ и μ , а также длина слова \widehat{w}_Σ , позволяющая отделить справочную часть K_0 от основной части $K_1(v)$. Затем производится изложенная в начале п. 6 обработка справочной части кода, результатом которой является (λ, μ) -представительное множество доопределений $\mathcal{D}_{\lambda, \mu}$, упорядоченное лексикографически.

На втором этапе сначала по основной части $K_1(v)$ кода находятся параметры d , b , z , t . Затем представляющие их слова \widehat{d} , \widehat{b} , \widehat{z} , \widehat{t} удаляются и остаток основной части разбивается на подслова длины $d + b + z$, задающие коды $K(v_i)$ кусков. Чтобы получить доопределение куска v_i , следует в соответствии с (27) разделить слово $K(v_i)$ на части длины d , b и z , найти в $\mathcal{D}_{\lambda, \mu}$ слово, номер которого при лексикографическом упорядочении задаётся первой частью кода, взять его начальное подслово, длина которого определяется второй частью, и осуществить в нём переименование символов в соответствии с третьей частью. Доопределение слова v образуется заменой кусков v_i полученными доопределениями кусков. ■

Приведём оценки характеристик кодирования при заданных параметрах λ и μ .

Лемма 8. Представленное кодирование частично определённых слов $v \in A^n$

1) обеспечивает длину кода

$$|K(v)| \leq h(v) + n \left(\frac{C_9 \log \lambda}{\mu - \log \lambda - 2} + \frac{\mu + C_9 \log \lambda}{\lambda} \right) + \lambda^{C_6} 2^\mu + 2 \log n;$$

2) выполнимо со сложностью, не превосходящей $(C_7)^\lambda + n\lambda^{C_{11}} 2^\mu$;

3) допускает декодирование, сложность которого не выше $(n + 2^\mu)\lambda^{C_{15}}$.

Доказательство. Утверждения 1 и 2 получены суммированием соответствующих результатов лемм 5 и 6. Член $C_{10}\lambda$, возникающий при суммировании длин, покрывается членом $\lambda^{C_6} 2^\mu$ (при подходящем выборе константы C_6) и потому в результат не включён.

Декодирование следует плану, описанному в лемме 7. На реализацию первого этапа затрачивается не более $\lambda^{C_{12}} 2^\mu$ операций (см. доказательство леммы 6). Возникающее (λ, μ) -представительное множество используется на втором этапе для декодирования кодов кусков $K(v_i)$. Выделение, обработка и декодирование кода $K(v_i)$ одного куска выполнимо с полиномиальной относительно λ сложностью. Число t кусков не превосходит n , поэтому декодирование всех кусков требует не более $n\lambda^{C_{15}}$ операций. Нетрудно видеть, что эта величина при подходящем выборе C_{15} покрывает также сложность других операций, используемых на втором этапе. Суммарная сложность обоих этапов составляет $\lambda^{C_{12}} 2^\mu + n\lambda^{C_{15}}$. Назначив константу C_{15} превосходящей C_{12} , приходим к утверждению 3 леммы 8. ■

Введём функцию $\phi(n) \rightarrow \infty$, в терминах которой будет сформулирована оценка средней длины кода (см. замечание 1). Имеет смысл рассматривать лишь функции $\phi(n)$ с порядком роста меньше $\left(\frac{\log n}{\log \log n}\right)^{1/2}$, ибо иначе остаточный член $O\left(\phi(n)\left(\frac{\log \log n}{\log n}\right)^{1/2}\right)$ не стремится к 0 и асимптотическая оптимальность кодирования не гарантируется. Будем считать, что функция $\phi(n)$ удовлетворяет условию

$$\phi(n) = o(\log^{1/2}(n)). \quad (34)$$

Назначим параметры λ и μ , положив

$$\lambda = \left\lfloor \frac{\log n}{\phi^2(n)} \right\rfloor; \quad (35)$$

$$\mu = \left\lfloor \frac{(\log n \log \log n)^{1/2}}{\phi(n)} \right\rfloor. \quad (36)$$

Запись $\lfloor x \rfloor$ означает ближайшее целое к x снизу. Используя в лемме 8 эти значения параметров, получаем следующий факт.

Лемма 9. Если растущая функция $\phi(n)$ удовлетворяет условию (34), то метод кодирования частично определённых слов длины n , использующий значения параметров λ и μ из (35), (36),

1) обеспечивает оценку длины кода

$$|K(v)| \leq h(v) + O\left(n\phi(n)\left(\frac{\log \log n}{\log n}\right)^{1/2}\right);$$

2) допускает кодирование и декодирование со сложностью $O(n^2)$.

Доказательство.

1) Перепишем утверждение 1 леммы 8 в виде

$$|K(v)| - h(v) \leq A_1 + A_2 + A_3 + A_4, \quad (37)$$

где $A_1 = \frac{C_9 n \log \lambda}{\mu - \log \lambda - 2}$; $A_2 = \frac{n(\mu + C_9 \log \lambda)}{\lambda}$; $A_3 = \lambda^{C_6} 2^\mu$; $A_4 = 2 \log n$.

Из (35) с учётом (34) следует, что $\lambda \rightarrow \infty$ и, кроме того, $\log \lambda \leq \log \log n$. Сравнение равенств (35) и (36) показывает, что $\mu \geq (\lambda \log \lambda)^{1/2} - 1$, а потому $\log \lambda = o(\mu)$ и имеют место асимптотические равенства $\mu - \log \lambda - 2 \sim \mu$ и $\mu + C_9 \log \lambda \sim \mu$. Принимая их во внимание, получаем

$$\begin{aligned} A_1 &\sim \frac{C_9 n \log \lambda}{\mu} \lesssim \frac{C_9 n \log \log n \phi(n)}{(\log n \log \log n)^{1/2}} \lesssim C_9 n \phi(n) \left(\frac{\log \log n}{\log n} \right)^{1/2}, \\ A_2 &\sim \frac{n\mu}{\lambda} \sim \frac{n(\log n \log \log n)^{1/2} \phi(n)}{\log n} \sim n \phi(n) \left(\frac{\log \log n}{\log n} \right)^{1/2}. \end{aligned}$$

Из последнего соотношения следует, в частности, что $\log A_2 \sim \log n$. В то же время $\log A_3 = C_6 \log \lambda + \mu \sim \mu = o(\log n)$, а потому $\log A_3 = o(\log A_2)$ и, тем более, $A_3 = o(A_2)$. Очевидно также, что $A_4 = o(A_2)$. Подстановка полученных соотношений в (37) даёт

$$|K(v)| - h(v) \lesssim (C_9 + 1)n\phi(n) \left(\frac{\log \log n}{\log n} \right)^{1/2},$$

а потому для некоторой константы C_{15} выполнено

$$|K(v)| \leq h(v) + C_{15} n \phi(n) \left(\frac{\log \log n}{\log n} \right)^{1/2}, \quad (38)$$

что влечет первое утверждение леммы.

2) Параметры λ и μ , заданные равенствами (35) и (36), удовлетворяют условиям $\lambda = o(\log n)$ и $\mu = o(\log n)$. Подстановки этих соотношений в оценки сложности кодирования и декодирования из п. 2 и 3 леммы 8 показывают, что каждая из этих оценок не превосходит $n^{1+o(1)} \leq O(n^2)$. Это даёт утверждение п. 2 леммы 9. ■

Следующее утверждение имеет место для недоопределённых источников $X = (A, P)$ в произвольном недоопределённом алфавите $A = \{a_T : T \in \mathcal{T}\}$.

Лемма 10. Если R — кодирование недоопределённого источника $X = (A, P)$ общего вида и для всякого слова $v \in A^n$ длина кодового слова удовлетворяет условию

$$K(v) \leq h(v) + G(n)$$

с некоторой функцией $G(n)$ (не зависящей от v), то для средней длины кода K справедлива оценка

$$\bar{l}_K^{(n)} \leq \mathcal{H}(P) + \frac{G(n)}{n}.$$

Доказательство. Обозначим через $p(P, \mathbf{r}, n)$ суммарную вероятность слов класса $\mathcal{K}_n(\mathbf{r})$, $\mathbf{r} = (r_T, T \in \mathcal{T})$. Вероятности $p(P, \mathbf{r}, n)$ имеют вид

$$p(P, \mathbf{r}, n) = \frac{n!}{\prod_{T \in \mathcal{T}} r_T!} \prod_{T \in \mathcal{T}} p_T^{r_T}$$

и образуют полиномиальное (мультиномиальное) распределение [24]. При каждом наборе вероятностей P оно обладает свойствами

$$\sum_{\mathbf{r}} p(P, \mathbf{r}, n) = 1, \quad \sum_{\mathbf{r}} p(P, \mathbf{r}, n) \frac{r_T}{n} = p_T \quad (T \in \mathcal{T}), \quad (39)$$

где r_T и p_T — компоненты наборов \mathbf{r} и P .

В выражении (3) для средней длины кода $\bar{l}_K^{(n)}$ сгруппируем слова $v \in A^n$, принадлежащие одному частотному классу, и, воспользовавшись условиями леммы и тем, что слова класса $\mathcal{K}_n(\mathbf{r})$ имеют одинаковую квазиэнтропию $h(v) = h_n(\mathbf{r}) = \mathbf{n}\mathcal{H}\left(\frac{\mathbf{r}}{n}\right)$, придём к оценке

$$\bar{l}_K^{(n)} \leq \sum_{\mathbf{r}} p(P, \mathbf{r}, n) \mathcal{H}\left(\frac{\mathbf{r}}{n}\right) + \sum_{\mathbf{r}} p(P, \mathbf{r}, n) \frac{G(n)}{n} = \Sigma_1 + \Sigma_2.$$

Оценим Σ_1 , применив к вогнутой функции \mathcal{H} (п. 2° леммы 1) неравенство Йенсена и второе равенство из (39). В результате имеем

$$\Sigma_1 \leq \mathcal{H}\left(\sum_{\mathbf{r}} p(P, \mathbf{r}, n) \frac{\mathbf{r}}{n}\right) = \mathcal{H}(P).$$

Первое равенство из (39) в применении к Σ_2 даёт $\Sigma_2 = \frac{G(n)}{n}$. Утверждение леммы получается заменой Σ_1 и Σ_2 в оценке для $\bar{l}_K^{(n)}$ величинами $\mathcal{H}(P)$ и $\frac{G(n)}{n}$. ■

Лемма 11. Для любого частично определённого источника $X = (A, P)$ и любой функции $\phi(n) \rightarrow \infty$ имеется способ кодирования K , обеспечивающий оценку средней длины кода

$$\bar{l}_K^{(n)} \leq \mathcal{H}(P) + O\left(\phi(n) \left(\frac{\log \log n}{\log n}\right)^{1/2}\right)$$

и допускающий кодирование и декодирование со сложностью $O(n^2)$.

Доказательство. Если $\phi(n)$ удовлетворяет условию (34), то нужное утверждение получается из леммы 9 путём применения леммы 10 к неравенству (38).

Если условие (34) нарушено, следует рассмотреть произвольную функцию $\phi'(n) \rightarrow \infty$, удовлетворяющую условиям $\phi'(n) \leq \phi(n)$ и (34). Утверждение леммы, справедливое по доказанному для функции $\phi'(n)$, будет верно и для исходной функции $\phi(n)$, которая не меньше $\phi'(n)$. ■

Утверждение теоремы 2 вытекает из леммы 11 при $\phi(n) = (\log \log n)^{1/2}$.

ЛИТЕРАТУРА

1. Галлагер Р. Теория информации и надёжная связь. М.: Сов. радио, 1974. 720 с.
2. Шоломов Л. А. Элементы теории недоопределённой информации // Прикладная дискретная математика. Приложение. 2009. №2. С. 18–42.
3. Шоломов Л. А. Сжатие частично определённой информации // Нелинейная динамика и управление. Вып. 4. М.: Физматлит, 2004. С. 377–396.
4. Ахо А., Хопкрофт Дж., Ульман Дж. Построение и анализ вычислительных алгоритмов. М.: Мир, 1979. 536 с.
5. Шоломов Л. А. Энтропия системы частично определённых последовательностей с вложенными областями определения // Нелинейная динамика и управление. Вып. 3. М.: Физматлит, 2003. С. 305–320.

6. Колмогоров А. Н. Три подхода к определению понятия «количество информации» // Проблемы передачи информации. 1965. Т. 1. № 1. С. 3–11.
7. Чашкин А. В. Дискретная математика. М.: Академия, 2012. 352 с.
8. Шеннон К. Э. Синтез двухполюсных переключательных схем // Работы по теории информации и кибернетике. М.: ИЛ, 1963. С. 59–101.
9. Лупанов О. В. Об одном подходе к синтезу схем — принципе локального кодирования // Проблемы кибернетики. Вып. 14. М.: Наука, 1965. С. 31–110.
10. Нечипорук Э. И. О сложности вентильных схем, реализующих булевские матрицы с неопределенными элементами // Докл. АН СССР. 1965. Т. 163. № 1. С. 40–43.
11. Шоломов Л. А. О реализации недоопределенных булевых функций схемами из функциональных элементов // Проблемы кибернетики. Вып. 21. М.: Наука, 1969. С. 215–226.
12. Андреев А. Е. О реализации частичных булевых функций схемами из функциональных элементов // Дискретная математика. 1989. Вып. 4. С. 36–45.
13. Чашкин А. В. Методы вычисления частичных булевых функций // Дискретные модели в теории управляющих систем: VII Междунар. конф. М.: МАКС Пресс, 2006. С. 390–404.
14. Шоломов Л. А. О функционалах, характеризующих сложность систем недоопределенных булевых функций // Проблемы кибернетики. Вып. 19. М.: Наука, 1967. С. 123–139.
15. Andreev A. E., Clementi A. E. F., and Rolim J. D. P. Worst-case hardness suffices for derandomization: A new method for hardness–randomness trade-offs // LNCS. 1997. V. 1256. 1997. P. 177–187.
16. Miltersen P. B. On the Shannon function for partially defined Boolean functions. <http://www.brics.dk/~bromille/Papers/index.html>. 1999.
17. Мадатьян Х. А. О реализации не всюду определенных k -значных матриц заданной «густоты» вентильными схемами глубины два // Методы дискретного анализа в теории булевых функций и схем. Вып. 35. Новосибирск: ИМ СО АН, 1980. С. 71–82.
18. Andreev A. E. Complexity of Nondeterministic Functions. BRICS report. Ser. RS-94-2. Febr. 1994. 47 p.
19. Чашкин А. В. Вычисление недоопределенных функций // Дискретная математика и ее приложения. Сб. лекций молодежных научных школ. Вып. VI. М.: ИПМ РАН, 2011. С. 29–40.
20. Krivevsky R. E. Occam's razor, partially specified Boolean functions, string matching, and independent sets // Information and Computation. 1994. Iss. 108. P. 158–174.
21. Кричевский Р. Е. Сжатие и поиск информации. М.: Радио и связь, 1989. 168 с.
22. Шоломов Л. А. Кодирование частично определенных дискретных источников без памяти // Докл. АН. 2004. Т. 397. № 2. С. 178–180.
23. Васильев Ю. Л., Глаголев В. В. Метрические свойства дизъюнктивных нормальных форм // Дискретная математика и математические вопросы кибернетики. Т. 1. М.: Наука, 1974. С. 99–148.
24. Крамер Г. Математические методы статистики. М.: Мир, 1975. 648 с.

REFERENCES

1. Gallager R. G. Information Theory and Reliable Communication. N.Y., Wiley Publ., 1968. 608 p.
2. Sholomov L. A. Elementy teorii nedoopredelennoy informatsii [Elements of the undetermined information theory]. Prikladnaya Diskretnaya Matematika. Prilozhenie, 2009, no. 2, pp. 18–42. (in Russian)

3. *Sholomov L. A.* Szhatie chastichno opredelennoy informatsii [Compression of partial defined information]. *Nelineinaya dinamika i upravlenie*, vol. 4, Moscow, Fizmatlit, 2004, pp. 377–396. (in Russian)
4. *Aho A., Hopcroft J., and Ulman J.* The Design and Analysis of Computer Algorithms. Addison-Wesley Publ. Co, 1976. 480 p.
5. *Sholomov L. A.* Entropiya sistemy chastichno opredelennih posledovatelnoctei s vlozhennimi oblastyami opredeleniya [Entropy of a system of partially defined sequences with nested domains]. *Nelineinaya Dinamika i Upravlenie*, vol. 3, Moscow, Fizmatlit Publ., 2003, pp. 305–320. (in Russian)
6. *Kolmogorov A. N.* Three approaches to the quantitative definition of information. *Problems Inform. Transmissions*, 1965, no. 1, pp. 1–7.
7. *Chashkin A. V.* Diskretnaya matematika [Discrete Mathematics]. Moscow, Academia, 2012. 352 p. (in Russian)
8. *Shannon C. E.* The synthesis of two-terminal switching circuits. *Bell Syst. Techn. J.*, 1949, vol. 28, no. 1. pp. 59–98.
9. *Lupanov O. B.* Ob odnom podhode k sintezu upravlyayushchikh sistem — printsipe lokalnogo kodirovaniya [On a certain approach to the synthesis of control systems — the principle of local coding]. *Problemy Kibernetiki*, iss. 14, Moscow, Nauka Publ., 1965, pp. 31–110. (in Russian)
10. *Nechiporuk E. I.* Complexity of gating circuits which are realized by Boolean matrices with undetermined elements. *Soviet Phis. Dokl.*, 1966, iss. 10, pp. 591–593.
11. *Sholomov L. A.* Realization of partial Boolean functions by circuits from functional elements. *Systems Theory Res.*, 1971, iss. 21, pp. 211–223.
12. *Andreev A. E.* On the complexity of the realization of Boolean functions by circuits of functional elements. *Discrete Math. Appl.*, 1991, vol. 1, iss. 3, pp. 251–261.
13. *Chashkin A. V.* Metody vychisleniya chastichnyh bulevykh funktsii [Methods for computing partial Boolean functions.] *Diskretnye Modeli v Teorii Upravlyayushchih Sistem: VII Mezhdunarodnaya Konferentsiya*. Moscow, MAKS Press, 2006, pp. 390–404. (in Russian)
14. *Sholomov L. A.* On functionals characterizing the complexity of systems of undetermined Boolean functions. *Systems Theory Res.*, 1970, iss. 20, pp. 123–140.
15. *Andreev A. E., Clementi A. E. F., and Rolim J. D. P.* Worst-case hardness suffices for derandomization: A new method for hardness–randomness trade-offs. *LNCS*, 1997, vol. 1256, pp. 177–187.
16. *Miltersen P. B.* On the Shannon function for partially defined Boolean functions. <http://www.brics.dk/~bromille/Papers/index.html>. 1999.
17. *Madatyan H. A.* O realizatsii ne vsyudu opredelennykh k -znachnykh matrits zadannoy “gustoty” ventilnymi shemami glubiny dva [On the implementation of not everywhere defined k -valued matrices of a given “density” by valve circuits of depth two]. *Metody Diskretnogo Analiza v Teorii Bulevykh Funktsiy i Skhem*, iss. 35, Novosibirsk, Institute of Mathematics Publ. House, 1980, pp. 71–82. (in Russian)
18. *Andreev A. E.* Complexity of nondeterministic functions. BRICS report Series, RS-94-2, Febr. 1994. 47 p.
19. *Chashkin A. V.* Vychislenie nedoopredelennykh funktsii [Computing of underdetermined functions.] *Diskretnaya Matematika i ee Prilozheniya. Sbornik Lektsiy Molodezhnykh Nauchnykh Shkol, VI*, Moscow, Keldysh Inst. of Appl. Math. Publ. House, 2011, pp. 29–40. (in Russian)
20. *Krichevsky R. E.* Occam’s razor, partially specified Boolean functions, string matching, and independent sets. *Information and Computation*, 1994, iss. 108, pp. 158–174.
21. *Krichevsky R.* Universal Compression and Retrieval. Kluwer Acad. Publ., 2010. 219 p.

22. *Sholomov L. A.* Encoding of partially defined discrete memoryless sources. *Doklady Mathematics*, 2004, vol. 70, no. 1, pp. 651–653.
23. *Vasil'ev Yu. L. and Glagolev V. V.* Metricheskie svoystva diz'yunktivnykh normal'nykh form [Metric properties of disjunctive normal forms]. *Discrete Mathematics and Mathematical Problems of Cybernetics*, vol. 1, Moscow, Nauka, 1974, pp. 99–148. (in Russian)
24. *Cramer H.* *Mathematical Methods of Statistics*. Princeton University Press, 1946. 575 p.