

УДК 512.54

А. А. Толстопятов¹

Медленный алгоритм кодирования и декодирования поля кратности при булевом сжатии файлов

Ключевые слова: булева алгебра, сжатие информации.

Предложен алгоритм, позволяющий для разложения заданного натурального числа на сумму фиксированного числа натуральных слагаемых построить лексикографически упорядоченную таблицу всех таких разложений и определить номер в этой таблице для любого разложения. Алгоритм работает медленно, так как приходится последовательно генерировать всю таблицу, но прост для программной реализации.

We suggest the algorithm that construct for representation of any natural number as a sum of fixed quantity of natural addends the lexicographically ordered table of all such representations; this algorithm also determines the number of such representation in this table. The algorithm works slowly, but is simple for program realization.

В работе [1] рассмотрен подход к сжатию дискретной информации без потерь, основанный на разбиении файла на кортежи длины n и объединении этих кортежей в буферы разной длины. При этом код каждого кортежа состоит из следующих трех полей:

- (1) *Поле принадлежности.* Это поле задается коэффициентами булева уравнения, решения которого дают кортежи, входящие в данный буфер.
- (2) *Поле кратности.* Это поле содержит информацию о том, сколько раз каждый кортеж входит в данный буфер. Для построения кода этого поля нужно определить номер разложения числа кортежей в буфере в сумму, содержащую столько слагаемых, сколько разных кортежей содержит данный буфер в лексикографически упорядоченной таблице всех таких разложений.
- (3) *Поле порядка.* Это поле содержит информацию о порядке расстановки кортежей внутри буфера. Его код задается номером в лексикографически упорядоченной таблице всех перестановок такого типа, который определяется полем кратности.

В настоящей работе рассматривается алгоритм, позволяющий построить код поля кратности. Если обозначить число кортежей в буфере через

¹Ивановский государственный университет; E-mail: khash2@mail.ru. Работа выполнена при финансовой поддержке РФФИ (проект 07-07-00155).

m , а число разных кортежей через s , то длина кода поля кратности L_2 , если задавать этот код номером в описанной ниже таблице, будет равна

$$L_2 = \log_2 D_{m-s}^s. \quad (1)$$

Как показано в [2], величина D_{m-s}^s может быть вычислена либо с использованием производящей функции

$$F(x) = \sum_{n=0}^{\infty} D_n^s x^n = \frac{1}{\prod_{k=1}^s (1-x^k)}, \quad (2)$$

$$D_n^s = \frac{1}{n!} \lim_{x \rightarrow \infty} \frac{d^n}{dx^n} \frac{1}{\prod_{k=1}^s (1-x^k)}, \quad (3)$$

либо с использованием рекуррентных формул:

1. $D_n^s = D_{n-s}^s + D_n^{s-1}, \quad s \leq n,$
2. $D_n^s = D_n^n, \quad s > n,$
3. $D_0^s = 1,$
4. $D_n^1 = 1,$
5. $D_n^2 = \left[\frac{n}{2} \right] + 1.$

(4)

Однако, в работе [2] были получены только асимптотические формулы для оценки L_2 , но не был построен алгоритм получения указанного выше номера. В настоящей работе этот пробел будет заполнен.

1. Постановка задачи

Пусть буфер файла разбит на m кортежей, причем среди них s различных. Назовем типом кратности n_1, n_2, \dots, n_s , где n_s — число вхождений в файл k -го кортежа, $k = 1, 2, \dots, s$. При этом $1 \leq s \leq m$ и выполнено соотношение

$$\sum_{k=1}^s n_k = m. \quad (5)$$

Если считать, что n_k лексикографически упорядочены, т. е.

$$n_k \leq n_{k+1}, \quad k = 1, 2, \dots, s-1,$$

то при фиксированных m и s все такие типы кортежей (или типы повторов) образуют таблицу:

$$\begin{aligned}
 & 1. \quad n_k = 1, \quad n_s = m - s + 1, \\
 & \quad \quad \quad \quad \quad \quad \quad k = 1, 2, \dots, s - 1, \\
 & 2. \quad n_k = 1, \quad n_{s-1} = 2, n_s = m - s \\
 & \quad \quad \quad \quad \quad \quad \quad k = 1, 2, \dots, s - 2, \\
 & \dots \\
 & D_{m-s}^2. \quad n_k = \left[\frac{m}{s} \right], \quad n_{p+l} = \left[\frac{m}{s} \right] + 1, \\
 & \quad \quad \quad \quad \quad \quad \quad k = 1, 2, \dots, p, \\
 & \quad \quad \quad \quad \quad \quad \quad p = \left(\left[\frac{m}{s} \right] + 1 \right) s - m, \\
 & \quad \quad \quad \quad \quad \quad \quad l = 1, 2, \dots, s - p.
 \end{aligned} \tag{6}$$

Задача построения кода поля кратности заключается в том, чтобы по заданному типу повторов n_1, n_2, \dots, n_s вычислить номер этого типа в таблице (6). Задача декодирования заключается в том, чтобы при заданных m и s по номеру в таблице (6) восстановить тип повторов n_1, n_2, \dots, n_s . В настоящей работе эти задачи будут решены с помощью алгоритма, генерирующего последовательно таблицу (6), и на каждом шаге сравнивающего заданный тип повтора n_1, n_2, \dots, n_s с соответствующей строкой таблицы (6).

2. Алгоритм кодирования

1. Обозначим через $\mathcal{N} = (n_1, n_2, \dots, n_s)$ тип повторов с заданным s . При этом m определяется из формулы (5).

2. Назовем $\tilde{\mathcal{N}} = (\tilde{n}_1, \tilde{n}_2, \dots, \tilde{n}_s)$ стандартной формой типа повторов, если $\tilde{\mathcal{N}}$ получается из \mathcal{N} такой перестановкой, чтобы выполнялось условие

$$\tilde{n}_1 \leq \tilde{n}_2 \leq \dots \leq \tilde{n}_s. \tag{7}$$

3. Определим последовательную процедуру генерации стандартных форм, положив, что переход от i -ой стандартной формы

$$\tilde{\mathcal{N}}_i = (\tilde{n}_1^i, \tilde{n}_2^i, \dots, \tilde{n}_s^i) \tag{8}$$

к $(i + 1)$ -ой стандартной форме

$$\tilde{\mathcal{N}}_{i+1} = (\tilde{n}_1^{i+1}, \tilde{n}_2^{i+1}, \dots, \tilde{n}_s^{i+1}) \tag{9}$$

проводится следующим образом.

3.1. Определяется число g первых a_j^i , $j = 1, 2, \dots, g$, которые не будут меняться при переходе от $\tilde{\mathcal{N}}_i$ к $\tilde{\mathcal{N}}_{i+1}$.

3.2. Проверяется первое условие

$$\tilde{n}_{s-1}^i - \tilde{n}_{s-t-1}^i > 1, \tag{10}$$

когда параметру t последовательно придаются значения $t = 0, 1, \dots, s - 2$.

3.3. Находится минимальное значение t ($\min t$) из $0, 1, 2, \dots, s-2$, при котором выполнено условие (10).

3.4. Тогда

$$g_1 = s - 2 - \min t, \quad t = 0, 1, 2, \dots, s - 2. \quad (11)$$

3.5. Проверяется второе условие

$$\tilde{n}_{s-1}^i - \tilde{n}_{s-t}^i = 2 \quad (12)$$

при последовательно изменяемых значениях $t = 2, 3, \dots, s - 1$.

3.6. Находятся минимальное значение t из $2, 3, \dots, s - 1$.

3.7. Тогда

$$g_2 = s - 1 - \min t, \quad t = 2, 3, \dots, s - 1. \quad (13)$$

3.8. Из g_1 и g_2 выбирается максимальное значение

$$g = \max(g_1, g_2). \quad (14)$$

3.9. Тогда алгоритм перехода от $\tilde{\mathcal{N}}_i$ из (8) к $\tilde{\mathcal{N}}_{i+1}$ из (9) имеет вид:

$$\begin{aligned} \tilde{n}_k^{i+1} &= \tilde{n}_k^i, \quad k = 1, 2, \dots, g; \\ \tilde{n}_k^{i+1} &= \tilde{n}_{g+1}^i, \quad k = g + 1, g + 2, \dots, s - 1; \\ \tilde{n}_s^{i+1} &= m - \sum_{k=1}^g \tilde{n}_k^i - (\tilde{n}_{g+1}^i + 1)(s - g - 1). \end{aligned} \quad (15)$$

4. Описание процедуры генерации стандартной формы $\tilde{\mathcal{N}}_{i+1}$ (9) из стандартной формы $\tilde{\mathcal{N}}_i$ (8) начинается с первой стандартной формы

$$\tilde{\mathcal{N}}_1 = (\tilde{n}_1^1, \tilde{n}_2^1, \dots, \tilde{n}_s^1), \quad (16)$$

где

$$\begin{aligned} \tilde{n}_k^i &= 1, \quad k = 1, 2, \dots, s - 1, \\ \tilde{n}_s^1 &= m - s + 1. \end{aligned} \quad (17)$$

5. Заканчивается процедура генерацией последней стандартной формы

$$\tilde{\mathcal{N}}_f = (\tilde{n}_1^f, \tilde{n}_2^f, \dots, \tilde{n}_s^f), \quad (18)$$

где

$$\begin{aligned} \tilde{n}_k^f &= \left[\frac{m}{s} \right], \quad k = 1, 2, \dots, g_1, \\ \tilde{n}_k^f &= 1 + \left[\frac{m}{s} \right], \quad k = g_1 + 1, \dots, s, \end{aligned} \quad (19)$$

где числа g_1 и g_2 определяются следующим образом:

$$g_1 = s - m + s \left[\frac{m}{s} \right] \quad \text{—} \quad (20)$$

число первых \tilde{n}_k^f , равных $\left[\frac{m}{s} \right]$.

$$g_2 = m - s \left[\frac{m}{s} \right] \quad \text{—}$$

число первых \tilde{n}_k^f равных $1 + \left[\frac{m}{s} \right]$.

6. На каждом шаге переходов от $\tilde{\mathcal{N}}_1$ к $\tilde{\mathcal{N}}_2$, от $\tilde{\mathcal{N}}_2$ к $\tilde{\mathcal{N}}_3$ и т. д. нужно сравнивать получаемую стандартную форму с исходной стандартной формой $\tilde{\mathcal{N}}$, номер которой нужно вычислить.

7. Если $\tilde{\mathcal{N}}_i$ не совпадает с $\tilde{\mathcal{N}}$, то нужно вычислить соответствующие p_i и z_i^j , где p_i — количество различных чисел в форме $\tilde{\mathcal{N}}_i$, а z_i^j — кратность каждого из p_i чисел в стандартной форме $\tilde{\mathcal{N}}_i$.

8. Тогда каждой стандартной форме $\tilde{\mathcal{N}}_i$ ставится в соответствие число

$$\sigma_i = \frac{s!}{p_i \prod_{j=1} z_i^j!}, \quad (21)$$

причем, если $\tilde{\mathcal{N}}_i$ совпадает с $\tilde{\mathcal{N}}$, то σ_i вычислять не нужно.

9. Зная все $\sigma_1, \sigma_2, \dots, \sigma_{i-1}$, найдем номер N_1 равный

$$N_1 = \sum_{k=1}^{i-1} \sigma_k. \quad (22)$$

10. Номер N типа повторов \mathcal{N} в полученной таблице будет равен

$$N = N_1 + N_2. \quad (23)$$

11. Для определения номера N_2 строится следующая процедура. Берется исходная форма $\mathcal{N} = (n_1, n_2, \dots, n_s)$.

12. Находится в $\mathcal{N} = (n_1, n_2, \dots, n_s)$ минимальное значение $n'_1, n'_2, \dots, n'_\nu$, где $n'_1 = n'_2, \dots, n'_\nu$.

13. Производится замена

$$n'_1 \rightarrow 1, n'_2 \rightarrow 1, \dots, n'_\nu \rightarrow 1. \quad (24)$$

14. Удаляются из $\mathcal{N} = (n_1, n_2, \dots, n_s)$ все $n'_1, n'_2, \dots, n'_\nu$.

15. В оставшемся множестве находятся минимумы значений $n_1^2, n_2^2, \dots, n_{s-\nu}^2$.

16. Производится замена

$$n_1^2 \rightarrow 2, n_2^2 \rightarrow 2, \dots, n_{s-\nu}^2 \rightarrow 2. \quad (25)$$

17. Удаляются из списка все $n_1^2, n_2^2, \dots, n_{s-\nu}^2$.

18. Процедура, описанная в пп. 12–17 повторяется до тех пор, пока $\mathcal{N} = (n_1, n_2, \dots, n_s)$ не превратится в пустое множество.

19. В $\mathcal{N} = (n_1, n_2, \dots, n_s)$ вместо n_k подставляются присвоенные им значения и получается перестановка j_1, j_2, \dots, j_s .

20. По известному алгоритму [3] по перестановке j_1, j_2, \dots, j_s находится ее номер N_2 .

21. Тогда, в соответствии с (23):

$$N = N_1 + N_2.$$

3. Алгоритм декодирования

Алгоритм декодирования должен по известным m, s, N восстановить тип повторов $\mathcal{N} = (n_1, n_2, \dots, n_s)$. По существу он уже содержится в алгоритме кодирования. Если сгенерированная при кодировании таблица (6) сохраняется, то достаточно просто найти в этой таблице строку с номером N . Она будет содержать искомый тип повторов $\mathcal{N} = (n_1, n_2, \dots, n_s)$. Тогда декодирование будет работать быстро, но потребует много памяти. При больших m этот вариант может оказаться невозможен технически. В этом случае можно при декодировании заново сгенерировать таблицу (6) согласно пп. 3–5 предыдущего раздела этой работы, закончив генерацию не последней стандартной формой $\widetilde{\mathcal{N}}_f$ из (18), (19), а формой с номером N . Тогда для декодирования никакой дополнительной памяти, кроме памяти для хранения номера N , не потребуется, но потребуется больше времени для декодирования, чем в первом варианте.

В заключение отметим, что тестирование рассмотренного в этой работе алгоритма показало, что при $m \sim 100$, когда среди $n_k, k = 1, 2, \dots, s$, есть много чисел, больших единицы, алгоритм работает недостаточно быстро, поэтому он и был назван медленным.

Список использованной литературы

1. Толстопятов А. А. О структуре дискретной информации и общих условиях ее сжатия // Вестник ИвГУ. – 2002. – Вып. 3. – С. 80–82.
2. Толстопятов А. А. Вычисление длины поля кратности при булевом сжатии файлов // Вестник ИвГУ. – 2004. – Вып. 3. – С. 71–76.
3. Толстопятов А. А. Быстрый алгоритм кодирования и декодирования поля порядка при булевом сжатии файлов // Математика и ее приложения: Журн. Иванов. матем. об-ва. – 2007. – Вып.1(4). – С. 35–46.

Поступила в редакцию 21.11.2007.