

УДК 621.327:681.5

В.В. Баранник¹, А.В. Хаханова²¹Харьковский университет Воздушных Сил им. И. Кожедуба²Харьковский национальный университет радиоэлектроники

ОЦЕНКА ИНФОРМАТИВНОСТИ ДВОИЧНЫХ МАССИВОВ НА ОСНОВЕ КОМБИНАТОРНОГО ПОДХОДА

Излагаются основные этапы методики оценки информативности двоичных данных на основе построения комбинаторной информационной модели двоичных массивов. Данная методика учитывает ограничения на число серий единиц в двоичных последовательностях плавающих размеров. Это позволило обосновать, что учет закономерностей по числу серий единиц позволяет сократить избыточность в условиях не стационарности статистических свойств исходных данных. Показано, что данный подход к оценке информативности двоичных данных обеспечивает сокращение избыточности в условиях: равновероятных символов алфавита; для наихудших условий выявленных структурных закономерностей.

Ключевые слова: число серий единиц, информационная модель, сжатие данных.

Введение

Постановка проблемы и анализ литературы. Существующие подходы относительно компактного представления двоичных данных, не адекватны современным требованиям процессов функционирования информационных систем [1, 2]. Увеличиваются объемы данных, которые предназначены для хранения и передачи в информационно-регистрирующих системах. Основной причиной является организация сжатия двоичных данных на основе устранения вероятностных видов избыточности. В тоже время количество такой избыточности зависит от типа законов распределения вероятностей появления символов в сообщении. Это приводит к снижению степени сжатия данных

Поэтому **актуальной научной задачей** является сокращение объемов двоичных данных без внесения погрешности.

Целью статьи является обоснование выбора подхода для построения технологии компрессии двоичных данных.

Разработка методики оценки информативности двоичных массивов

Для дополнительного увеличения коэффициента сжатия, обрабатываемых двоичных данных предлагается дополнительно проводить сокращение избыточности в двоичных матрицах G , $G = \{g_{k\ell}\}$, $k = \overline{1, n}$, $\ell = \overline{1, n}$, $g_{k\ell} \in \{0; 1\}$, где $g_{k\ell}$ – $(k; \ell)$ -й элемент матрицы знаков. Двоичные матрицы имеют низкую избыточность $R_{ст}$, определяемую на основе учета статистических закономерностей $R_{ст} = (N_0 - N_{ст})/N_0 \rightarrow 0$, где N_0 – количество информации, в среднем приходящееся на один элемент матрицы, для случая, когда закономерности не выявлены, $N_0 = 1$ бит; $N_{ст}$ – среднее количество

информации, содержащееся в одном элементе с учетом ограничений на закон распределения вероятностей $p_{k\ell}$ появления элементов $g_{k\ell}$: $N_{ст} = p(0)\log_2 p(0) + p(1)\log_2 p(1) \rightarrow N_0$, где $p(0)$, $p(1)$ – вероятности появления соответственно нулевого и единичного элементов. Это снижает потенциальные возможности относительно обеспечения необходимой степени компрессии двоичных матриц. Следовательно, сжатие двоичных матриц достигается за счет сокращения комбинаторной избыточности, независимой от статистических свойств матриц G .

Анализируемая характеристика двоичных матриц должна учитывать взаимное расположение нулевых и единичных элементов и их количество. Причем количество нулевых элементов может равняться количеству единичных. Количество информации N в двоичных последовательностях с учетом анализируемых ограничений должно удовлетворять следующим требованиям:

1) количество информации на один элемент должно быть меньше единицы для источника информации с равномерным появлением единиц и нулей

$$N \ll 1 \text{ для } p(0) \approx p(1); \quad (1)$$

2) величина N для произвольных значений анализируемой характеристики η должна быть меньше единицы

$$\max_{\eta_{\min} \leq \eta \leq \eta_{\max}} \{N\} < 1, \quad (2)$$

где η_{\min} , η_{\max} – соответственно минимальное и максимальное значения величины η .

Предлагается рассматривать выявления закономерностей на основе характеристики количества серий η . Признак серии задается системой выражений:

$$\begin{cases} g_{k-1, \ell} \neq g_{k\ell} \rightarrow k \leq n; \\ g_{k, \ell} \neq g_{1, \ell+1} \rightarrow k > n. \end{cases}$$

Вычисление количества серий единиц η для v элементов матрицы G :

- на нулевом шаге $g_{0,\ell} = 0$; $\eta_\ell = 0$;
- на k -м шаге значения числа серий увеличивается на 1, $\eta_{k\ell} = \eta_{k-1,\ell} + 1$, если $g_{k-1,\ell} < g_{k\ell}$;
- в противном случае $\eta_{k\ell} = \eta_{k-1,\ell} + 0$, если $g_{k-1,\ell} \geq g_{k\ell}$;
- для конечного шага при $k=v$ получаем искомого значение количества серий единиц $\eta = \eta_v$ для j -го столбца.

Количество $V_{v,\eta}$ двоичных последовательностей G_v , содержащих v элементов и η серийных перепадов равно [3, 4]:

$$V_{v,\eta} = (v+1)! / (2\eta)! (v+1-2\eta)! \quad (3)$$

В соответствии с (3) количество информации H в среднем приходящееся на один элемент двоичной последовательности с η сериями единиц равно

$$H = \left(\sum_{\xi=1}^{v+1} \log_2 \xi - \sum_{\xi=1}^{2\eta} \log_2 \xi - \sum_{\xi=1}^{v+1-2\eta} \log_2 \xi \right) / v \quad (4)$$

Рассмотрим соответствие H условиям (1) и (2):
 1) для фиксированного v максимальное значение $V_{v,\eta}$ двоичных последовательностей достигается для $\eta = \eta_{cp} = [(v+1)/4]$. Это вытекает из того, что величина $V_{v,\eta}$ для $0 \leq \eta \leq [(v+1)/2]$ имеет один максимум в области среднего значения η . Поэтому для проверки условия (2) необходимо определить значение величины $V_{v,\eta_{cp}}$:

$$V_{v,\eta_{cp}} = \begin{cases} (v+1)! / (v/2)! (v/2+1)!, \\ \rightarrow v \bmod(2) = 0; \\ (v+1)! / ((v+1)/2)! ((v+1)/2)!, \\ \rightarrow v \bmod(2) \geq 1, \end{cases} \quad (5)$$

где $v \bmod(2) = 0$, $v \bmod(2) \geq 1$ – условия, когда величина v принимает соответственно четное и нечетные значения.

Значит, для $\eta = \eta_{cp} = [(v+1)/4]$ величина H будет принимать максимальное значение, равное:

- если $v \bmod(2) = 0$, то

$$H_{v,\eta_{cp}} = \left(\sum_{\xi=1}^{v+1} \log_2 \xi - 2 \sum_{\xi=1}^{v/2} \log_2 \xi - \log_2 (v/2+1) \right) / v;$$

- если $v \bmod(2) \geq 1$, то

$$H_{v,\eta_{cp}} = \left(\sum_{\xi=1}^{v+1} \log_2 \xi - 2 \sum_{\xi=1}^{(v+1)/2} \log_2 \xi \right) / v, \quad (6)$$

где $H_{v,\eta_{cp}}$ – количество информации, приходящееся на один элемент матрицы в случае когда, число серий единиц равно $\eta = \eta_{cp} = [(v+1)/4]$.

Условие (2) выполняется, если

$$H_{v,\eta_{cp}} < 1. \quad (7)$$

Правая часть неравенства (7) соответствует варианту, когда количество допустимых двоичных последовательностей равно 2^v . Поэтому неравенство (7) будет выполняться для

$$V_{v,\eta_{cp}} < 2^v. \quad (8)$$

Рассмотрим сумму величин $V_{v,\eta}$ по всему диапазону значений числа серий единиц $0 \leq \eta \leq [(v+1)/2]$:

1. Если v четное, то $\eta = \overline{0, v/2}$. Введем вспомогательные переменные $\alpha = (v+1)$ и $\beta = 2\eta$, то

$$\sum_{\beta=0}^{(\alpha-1)} (\alpha)! / ((\beta)! (\alpha-\beta)!). \quad (9)$$

Переменная β принимает все значения в диапазоне $\beta = \overline{0, (\alpha-1)}$. В то время как величина 2η принимает в этом диапазоне только положительные четные значения. Поэтому

$$\sum_{\eta=0}^{v/2} V_{v,\eta} = \sum_{\beta=0}^{(\alpha-1)} h_{\alpha,\beta}; \quad (10)$$

$$h_{\alpha,\beta} = \begin{cases} (\alpha)! / ((\beta)! (\alpha-\beta)!), & \rightarrow (-1)^\beta = 1; \\ 0, & \rightarrow (-1)^\beta = -1. \end{cases}$$

Тогда с учетом свойства сочетаний без повторов

$$\sum_{\beta=0}^{\alpha} ((-1)^\beta (\alpha)! / ((\beta)! (\alpha-\beta)!)) = 0 \quad (11)$$

получим $\sum_{\eta=0}^{v/2} V_{v,\eta} = 2^{\alpha-1} = 2^v$. В тоже время

$V_{v,\eta_{cp}} < \sum_{\eta=0}^{v/2} V_{v,\eta}$. Отсюда следует, что для четного

v неравенство (8) выполняется.

2. Если v нечетное, то $\eta = \overline{0, (v+1)/2}$, а

$$\sum_{\eta=0}^{(v+1)/2} V_{v,\eta} = \sum_{\eta=0}^{(v+1)/2} \frac{(v+1)!}{(2\eta)! (v+1-2\eta)!}.$$

С учетом свойств, заданных соотношениями (10) и (11), получим

$$\sum_{\eta=0}^{v/2} V_{v,\eta} = 2^{\alpha-1} = 2^v.$$

Следовательно, неравенство (8) выполняется для нечетного v . Таким образом, условие (2) для предлагаемой количественной характеристики «число серий единиц» выполняется.

Рассмотрим условие (1) применительно к характеристике «число серий единиц». Количество $V(s)_{v,\eta}$ двоичных последовательностей для числа единиц, равного s , распределенных по η сериям равно

$$V(s)_{v,\eta} = \frac{(s-1)!}{(\eta-1)!(s-\eta)!} \times \frac{(v-s+1)!}{(\eta)!(v-s+1-\eta)!},$$

$$\eta = 0, \min(s; v/2). (12)$$

В случае равных вероятностей $p(0) \approx p(1)$ имеем одинаковое количество нулевых и единичных элементов, равное $v/2$. Для такого условия выражение (12) будет иметь вид:

1. Если v – четное ($\eta = 0, s$), то

$$V(v/2)_{v,\eta} = \frac{((v/2)-1)!}{(\eta-1)!((v/2)-\eta)!} \times \frac{((v/2)+1)!}{(\eta)!((v/2)+1-\eta)!}. (13)$$

Поскольку сомножители в правой части выражения (13) представляют собой сочетания без повторений, то они достигают максимума при $\eta/2$. Значит максимальное значение величина $V(s)_{v,\eta}$ при фиксированных значениях v и s достигает для $\eta = s/2 = v/4$:

$$V(v/2)_{v,v/4} = \frac{((v/2)-1)!}{((v/4)-1)!((v/4))!} \times \frac{((v/2)+1)!}{((v/4))!((v/4)+1)!}.$$

В соответствии с данным выражением количества информации $\bar{H}(v/2)_{v,v/4}$, приходящееся на один элемент двоичной последовательности, содержащей $\eta = v/4$ серий единиц и $v/2$ единичных элементов

$$\bar{H}(v/2)_{v,v/4} = \log_2 \left(\frac{((v/2)-1)!}{((v/4)-1)!((v/4))!} \times \frac{((v/2)+1)!}{((v/4))!((v/4)+1)!} \right) / v.$$

Поскольку количество η серий единиц дополнительно ограничено величиной s , то выполняется

$$\sum_{\eta=1}^s V(s)_{v,\eta} = \sum_{\eta=1}^{v/2} V(v/2)_{v,\eta} < \sum_{\eta=0}^{v/2} V_{v,\eta}. (14)$$

Правая часть неравенства (14) имеет ограничение $\sum_{\eta=0}^{v/2} V_{v,\eta} = 2^v$. Значит $\sum_{\eta=1}^s V(s)_{v,\eta} < 2^v$, т.е.

$$2 \sum_{\xi=1}^{((v/2)-1)} \log_2 \xi - \sum_{\xi=(v/2)}^{((v/2)+1)} \log_2 \xi - 2 \sum_{\xi=1}^{((v/4)-1)} \log_2 \xi - \sum_{\xi=(v/4)}^{((v/4)+1)} \log_2 \xi - 2 \sum_{\xi=1}^{(v/4)} \log_2 \xi < v.$$

ОЦІНКА ІНФОРМАТИВНОСТІ ДВІЙКОВИХ МАСИВІВ НА ОСНОВІ КОМБІНАТОРНОГО ПІДХОДУ

В.В. Баранник, Г.В. Хаханова

Висловлюються основні етапи методики оцінки інформативності двійкових даних на основі побудови комбінаторної інформаційної моделі двійкових масивів. Дана методика враховує обмеження на число серий одиниць в двійкових послідовностях плаваючих розмірів. Це дозволило обґрунтувати, що облік закономірностей за числом серий одиниць дозволяє скоротити надмірність в умовах не стаціонарності статистичних властивостей початкових даних. Показано, що даний підхід до оцінки інформативності двійкових даних забезпечує скорочення надмірності в умовах: рівноімовірних символів алфавіту; для якнайгірших умов виявлених структурних закономірностей.

Ключові слова: число серий одиниць, інформаційна модель, стиснення даних.

ESTIMATION OF INFORMING OF BINARY ARRAYS ON THE BASIS OF COMBINATORY APPROACH

V.V. Barannik, A.V. Hahanova

The basic stages of method of estimation of informing of binary information are expounded on the basis of construction of combinatory informative model of binary arrays. This method takes into account limit on the number of cerouss of units in binary sequences of floating sizes. It allowed to ground, that the account of conformities to the law on the number of cerouss of units allows to shorten surplus in the conditions of not stationarity of statistical properties of basic data. It is rotined that this fitting to the estimations of informing of binary information are provided by reduction of surplus in terms: equal probabilities characters of alphabet; for the worst terms of the exposed structural conformities to the law.

Keywords: number of series of units, informative model, compression of data.

Следовательно, максимальное значение количества информации $\bar{H}(v/2)_{v,v/4}$ будет меньшим единицы. Условие (1) выполняется. По аналогии доказывается, что условие (1) выполняется для нечетного v . Таким образом, двоичные матрицы данных имеют комбинаторную избыточность, обусловленную ограничениями на число серий единиц.

Выводы

Построена комбинаторная модель оценки информативности двоичных матриц данных. Обосновано, что учет закономерностей по числу серий единиц позволяет сократить избыточность в условиях не стационарности статистических свойств исходных фрагментов данных.

Научная новизна состоит в том, что получила дальнейшее развитие методика оценки информативности двоичных данных на основе построения комбинаторной информационной модели двоичных массивов. В отличие от других моделей учитываются ограничения на число серий единиц в двоичных последовательностях плавающих размеров. Это позволило обосновать, что учет закономерностей по числу серий единиц позволяет сократить избыточность в условиях не стационарности статистических свойств исходных данных.

Список литературы

1. Прэтт У. Цифровая обработка изображений: Том 1, 2. – М.: Мир, 1985. – 736 с.
2. Методы сжатия данных. Устройство архиваторов, сжатие изображений и видео / В. Ватолин, А. Ратушняк и др. – М.: ДИАЛОГ – МИФИ, 2002. – 384 с.
3. Королев А.В., Баранник В.В. Оценка количества информации изображения по числу серий одинаковых элементов // Системы обработки информации. – Х.: НАНУ, ПАНМ, ХВУ, 2002. – Вып. 2 (18). – С. 43-46.
4. Баранник В.В., Юдин А.К. Рекуррентное двухпризнаковое двоичное полиадическое кодирование // Открытые информационные и компьютерные интегрированные технологии. – Х.: НАКУ «ХАИ», 2006. – Вып. 33. – С. 22-28.

Поступила в редколлегию 12.08.2008

Рецензент: д-р техн. наук, проф. Ю.В. Стасев, Харьковский университет Воздушных Сил им. И. Кожедуба, Харьков.