



ANÁLISIS TEÓRICO DEL ALGORITMO DE COMPRESIÓN LLRUN PARA LA COMPRESIÓN DE CADENAS DE BITS SPARSE

(Theoretical analysis LLRUN compression algorithm for compressing bit strings
SPARSE)

Recibido: 14/04/2015 Aceptado: 05/06/2014

Rincón, Carlos

Universidad del Zulia, Venezuela
crincon@fec.luz.edu.ve

Bracho, David

Universidad del Zulia, Venezuela
drbracho@fec.luz.edu.ve

Acurero, Alfredo

Universidad del Zulia, Venezuela
aacurero@fec.luz.edu.ve

RESUMEN

La presente investigación tuvo como finalidad realizar un análisis teórico sobre el funcionamiento del algoritmo de compresión LLRUN, para determinar su posible uso en la compresión de cadenas de bits sparse. Fraenkel y Klein (1985) plantearon un método para comprimir cadenas de bits sparse mezclando las técnicas run length encoding (RLE), Gamma coding y Huffman coding, mostrando un rendimiento prometedor. El resultado del presente estudio permitió determinar que LLRUN ofrece una relación de compresión teórica entre 87,5% y 50%, siendo una herramienta efectiva para la compresión de cadenas de bits sparse.

Palabras clave: compresión, cadenas de bits sparse, LLRUN, rendimiento.

ABSTRACT

The purpose of this research was to make a theoretical analysis about the performance of LLRUN compression algorithm on sparse bit strings. Fraenkel and Klein (1985) proposed a method to compress sparse bit strings, using run length encoding (RLE), Elias - gamma coding and Huffman coding, with promising performance. The results of this research showed that LLRUN offers a theoretical compression ratio between 87,5% and 50%, being an effective tool to compress sparse bit strings.

Keywords: compression, sparse bit strings, LLRUN, performance.



INTRODUCCI  N

En la actualidad la compresi  n de datos es una de las   reas de las ciencias computacionales de mayor auge, como consecuencia de la necesidad creciente de almacenamiento / transmisi  n de informaci  n por parte de los usuarios.

Un caso particular de la compresi  n de datos es la compresi  n de cadenas de bits sparse.

Fraenkel y Klein (1985) plantearon el m  todo LLRUN para comprimir cadenas de bits sparse, mezclando las t  cnicas Run Length Encoding (RLE), Elias - Gamma Coding y Huffman Coding.

A pesar de mostrar resultados prometedores, es escasa la informaci  n de   ndole acad  mica sobre el funcionamiento de este algoritmo, siendo esta una de las razones que fundament   la presente investigaci  n.

Rinc  n y col. (2009) presentaron un nuevo algoritmo de compresi  n basado en la teor  a de la informaci  n. Estudios adicionales sobre el rendimiento de este algoritmo, realizado por Rinc  n y col. (2012), permitieron determinar un efecto adverso del tama  o del alfabeto en el rendimiento objeto de estudio, generando una cadena de bits sparse como resultado de la ejecuci  n del algoritmo.

La necesidad de contar con una herramienta para la compresi  n de cadenas de bits sparse fue la otra raz  n que llev   a los investigadores a estudiar el algoritmo LLRUN.

TEOR  A DE LA INFORMACI  N

El art  culo "Una teor  a matem  tica de la comunicaci  n" de Shannon (1948) revolucion   el concepto que hasta la fecha se ten  a sobre la informaci  n.

Sin embargo, en la actualidad, la definici  n que se maneja sobre la informaci  n (inclusive por miembros de la comunidad de la computaci  n) es la uni  n de un conjunto de datos organizados que representan algo. El problema b  sico de esta definici  n es que no estipula el mecanismo que permita "medir" la informaci  n.

A pesar de la complejidad matem  tica impl  cita en los estudios de Shannon, la fundamentaci  n de estos es relativamente sencilla. Para Shannon, la informaci  n puede ser medida mediante una funci  n representada por la inversa de la probabilidad de ocurrencia del evento que se est   estudiando. Espec  ficamente, Shannon define la informaci  n mediante la f  rmula $I(x) = \log_2(1/P(x))$, donde x es el evento estudiado.

COMPRESI  N DE DATOS

La compresi  n de datos es una t  cnica que permite representar la informaci  n generada por una fuente de datos con una menor cantidad de bits, utilizando para este fin, c  digos   ptimos.



Las soluciones de compresi n se dividen en 2 pasos:

(a) Algoritmo de compresi n: es aquel que determina un c digo  ptimo que permita representar los datos con la menor cantidad de bits y que genera un archivo comprimido bas ndose en los c digos  ptimos encontrados.

(b) Algoritmo de descompresi n: es aquel que se encarga de transformar el archivo comprimido en la informaci n original.

RELACI N DE COMPRESI N

Es una m trica que permite determinar el rendimiento de un algoritmo de compresi n. Matem ticamente puede definirse de dos formas:

(a) Relaci n n a 1: $RC = TO/TC$ y

(b) Como porcentaje: $RC = (TO-TC)/TO$, donde TO = tama o del archivo original y TC = tama o del archivo comprimido.

A pesar que la relaci n de compresi n es planteada como una funci n cuyas variables independientes son el tama o del archivo original y el tama o del archivo comprimido, esta m trica puede ser simplificada considerando que el tama o del archivo original es igual al n mero de s mbolos del archivo multiplicado por el tama o promedio de los s mbolos sin comprimir (N_{SC}), y que el tama o del archivo comprimido es igual al n mero de s mbolos del archivo multiplicado por el tama o promedio de los s mbolos comprimidos (N_C). Considerando lo antes expuesto, se puede concluir que $RC = N_{SC}/N_C$ (relaci n n a 1) y $RC = (N_{SC} - N_C)/N_{SC} * 100$ (porcentaje).

RUN LENGTH ENCODING (RLE)

Golomb (1966) plante  un m todo para codificar una fuente donde se presenten una cantidad sucesiva de eventos desfavorables. Al aplicar la teor a de la informaci n, el algoritmo RLE codifica la cantidad de eventos desfavorables (Run Lengths), minimizando la informaci n generada.

El comportamiento probabil stico de la codificaci n de un Run Length de tama o n es igual p^nq , para $n = 0,1,2,3, \dots$, siendo similar a la distribuci n geom trica.

Para la presente investigaci n, se utiliz  la versi n RLE de 8bit, dividiendo la cadena sparse en octetos y codificando la posici n dentro de cada octeto de los bits encendidos.

Es importante resaltar que al dividir la cadena en octetos, se minimiza a 3 la cantidad de bits necesaria para representar una posici n, pero se tiene ambigüedad (c digo no es de traducci n  nica) en el caso que la posici n del primer bit encendido del octeto actual es mayor que la  ltima posici n del  ltimo bit encendido del octeto anterior. Ejemplo:



Cadena de Bits:

01000000 00000001 00010000,

Cadena RLE 8bits= 2 8 4

Decodificaci n cadena RLE =

01000001 00010000

Para evitar este error, se decidi  sumar ocho (8) a todas las posiciones encendidas de un octeto, s lo en el caso que la posici n del primer bit encendido del octeto actual sea mayor que la  ltima posici n del  ltimo bit encendido del octeto anterior. Al conseguir el primer valor mayor que 8, se conoce que es el nuevo octeto aquel cuyos bits encendidos ser n aquellos valores de las posiciones (sucesivas a la primera encontrada) que sean mayores que 8. La posici n real es aquella que se obtiene de restar 8 a las posiciones mayores que 8.

Al aplicar esta modificaci n al ejemplo anterior tenemos:

Cadena de Bits:

01000000 00000001 00010000,

Cadena RLE 8bits= 2 16 4

Decodificaci n cadena RLE =

01000000 00000001 00010000

CODIFICACI N ELIAS – GAMMA

Elias (1975), plante  un m todo para representar n meros enteros utilizando $2\log_2(X) + 1$ bits, siendo X el n mero entero a codificar.

En resumen, el m todo de codificaci n consiste en representar el n mero entero en 2 partes.

La primera parte (el prefijo) representa la potencia de 2 m s cercana al n mero X, tal que $2^N \leq X$. La segunda parte (el sufijo), es la representaci n binaria del n mero que se debe sumar a 2^N para ser igual a X. El tama o del sufijo es igual a la cantidad de bits del prefijo - 1.

Para la implementaci n del algoritmo LLRUN en la presente investigaci n, se utiliz  la representaci n en codificaci n Elias – Gamma, de los primeros 16 n meros (ver Tabla 1).

Tabla 1. Codificaci n Elias - Gamma

Prefijo	Sufijo	C�culo	Valor
1	N/A	$2^0 + N/A$	1
01	0	$2^1 + 0$	2
01	1	$2^1 + 1$	3
001	00	$2^2 + 0$	4
001	01	$2^2 + 1$	5
001	10	$2^2 + 2$	6
001	11	$2^2 + 3$	7
0001	000	$2^3 + 0$	8
0001	001	$2^3 + 1$	9
0001	010	$2^3 + 2$	10
0001	011	$2^3 + 3$	11
0001	100	$2^3 + 4$	12
0001	101	$2^3 + 5$	13
0001	110	$2^3 + 6$	14
0001	111	$2^3 + 7$	15
00001	0000	$2^4 + 0$	16

Fuente: elaboraci n propia.

Para los prefijos, se utiliza una representaci n inversa de la cadena de bits que representa la potencia de 2, con la finalidad de utilizar el bit encendido como delimitador entre el prefijo y el sufijo.

En la implementaci n de la codificaci n Elias – Gamma de la presente investigaci n, se descart  el sufijo del valor 16, por ser el  nico valor cuyo prefijo tiene 5 bits y ser igual a cero.

ALGORITMO DE COMPRESI N DE HUFFMAN

Huffman (1952) plante  una t cnica que permite generar c digos  ptimos para realizar el proceso de compresi n. Esta t cnica de compresi n sin p rdida se conoce en la actualidad como el algoritmo de compresi n de Huffman, el cual mantiene su vigencia a pesar de la introducci n de nuevos algoritmos para generar c digos  ptimos.

Los pasos que estipul  Huffman para obtener c digos  ptimos son:

1. Calcular las frecuencias de los s mbolos que aparecen en el mensaje.
2. Ordenar los s mbolos seg n su frecuencia de menor a mayor. Deben considerarse criterios para determinar el orden en los casos donde los s mbolos tienen la misma frecuencia.



3. Crear un s mbolo compuesto que estar  conformado por los s mbolos con menor frecuencia de la lista producto del paso 2. La frecuencia de este s mbolo compuesto ser  igual a la suma de las frecuencias de los s mbolos que lo conforman.

4. Reemplazar los s mbolos seleccionados por el s mbolo compuesto y reordenar los s mbolos del mensaje seg n su frecuencia.

5. Ir al paso 3 hasta que se tenga un  rbol binario cuyo nodo padre sea un s mbolo compuesto con una frecuencia igual al n mero de s mbolos del mensaje.

6. Etiquetar las ramas del  rbol con el valor de 0 y 1 (se debe establecer y mantener el criterio sobre qu  valor que va a tomar la rama izquierda y el valor que va a tomar la rama derecha).

7. Construir los c digos de Huffman (1952) de los s mbolos del mensaje uniendo las etiquetas que conforman el camino desde la ra z del  rbol hasta el la hoja que representa dicho s mbolo. El  rbol resultante tendr  tantos nodos hoja como s mbolos tenga el mensaje.

ALGORITMO LLRUN

Fraenkel y Klein (1985) presentaron un nuevo m todo para la compresi n de cadenas de bits sparse.

Los pasos de este algoritmo son los siguientes:

1. Aplicar RLE a la cadena de bits sparse, utilizando compresi n  ndice para descartar los octetos sin bits encendidos.

2. Aplicar la codificaci n Elias – Gamma a la cadena RLE.

3. Aplicar el algoritmo de Huffman (1952) a los prefijos producto de la codificaci n Elias - Gamma.

En la siguiente secci n, se presenta el funcionamiento del algoritmo.

RENDIMIENTO DEL ALGORITMO LLRUN

El rendimiento de un algoritmo de compresi n se mide determinando su mejor y peor caso. El mejor caso se considera cuando en el mensaje s lo aparece un s mbolo del alfabeto mientras que el peor caso es cuando aparecen n s mbolos de manera equiprobable en el mensaje.

Mejor Caso:

Para el algoritmo LLRUN, considerando que se trata de una t cnica de compresi n para cadenas de bits sparse, se utiliz  una cadena de 32 bits donde cada uno de los 4 octetos tiene encendido s lo el bit m s significativo:



Cadena de bits:

10000000 10000000 10000000 10000000

Al ejecutar el RLE de 8bits se obtiene el siguiente arreglo:

Arreglo RLE: 1 1 1 1

Al aplicar la codificaci n Elias - Gamma al arreglo RLE se obtiene la siguiente cadena:

Cadena GAMMA: 1 1 1 1 (4 bits)

Al aplicar el algoritmo de Huffman (1952) a los prefijos de la cadena Gamma se obtiene el siguiente c digo:

C (1) = 0

Al sustituir el c digo producto de la ejecuci n de Huffman (1952) en la cadena GAMMA se obtiene la cadena:

Cadena Comprimida = 0 0 0 0 (4bits)

Es importante resaltar que para este caso ninguno de los c digos GAMMA tiene sufijo dado a que el n mero 1 se representa como $2^0 + 0$.

La relaci n de compresi n para el mejor caso es:

RC (mejor caso) = $(32 \text{ bit} - 4 \text{ bits}) / 32 \text{ bit} * 100 = 87,5 \%$

Peor Caso:

Para el peor caso, se consider  una cadena de 64 bits donde cada uno de los octetos tiene encendido un bit diferente. La raz n por la que se consideraron los 64 consisti  en cumplir la caracter stica de equiprobabilidad que deben cumplir los s mbolos en el mensaje.

Cadena de bits: 10000000 01000000 00100000 00010000 00001000 00000100 00000010 00000001

Al ejecutar el RLE de 8bits se obtiene el siguiente arreglo:

Arreglo RLE: 1 2 3 4 5 6 7 8, sin embargo dada la implementaci n realizada del algoritmo RLE de 8bit que permite diferenciar el octeto al que pertenece una posici n encendida, el arreglo final es:

Arreglo RLE: 1 10 3 12 5 14 7 16

Al aplicar la codificaci n Elias - Gamma al arreglo RLE se obtiene la siguiente cadena:



Cadena GAMMA: 1 0001 010 01 1 0001 100 001 01 0001 110 001 11 00001

(40 bits)

Al aplicar el algoritmo de Huffman (1952) a los prefijos de la cadena Gamma se obtienen los siguientes c odigos:

$C(0001) = 11$, $C(001) = 01$, $C(01) = 101$, $C(00001)=00$, $C(1)=100$

Al sustituir los c odigos producto de la ejecuci on de Huffman (1952) en la cadena GAMMA se obtiene la cadena:

Cadena Comprimida = 100 11 010 101 1 11 100 01 01 11 110 01 11 00 (32 bits)

La relaci on de compresi on para el peor caso es:

RC (peor caso) = $(64 \text{ bit} - 32 \text{ bits})/32 \text{ bit} * 100 = 50 \%$

Es importante resaltar que para el c alculo de la relaci on de compresi on no se considera la informaci on adicional que necesita Huffman (1952) y la codificaci on por  ndice para aplicar la descompresi on. Las razones que fundamentan este criterio se derivan del aporte m ınimo que significa esta informaci on adicional, cuando se comprimen archivos de Megabits o Gigabits de informaci on.

CONCLUSIONES Y TRABAJOS FUTUROS

Luego del an alisis te orico realizado, se determin o que el rendimiento que ofrece el algoritmo LLRUN se encuentra entre el 87,5% y el 50% (sin considerar la informaci on de cabecera de las t ecnicas de compresi on utilizadas). Este rendimiento permite proponer a LLRUN como una opci on efectiva para la compresi on de cadenas de bits sparse.

Seg un Rinc on y col. (2012), el algoritmo de Compresi on Probabil stico, basado en la posici on de los s ımbolos, genera cadenas de bits sparse para alfabetos con tama o mayor a 15. La pr oxima investigaci on consistir a en aplicar el algoritmo LLRUN para tratar de optimizar el rendimiento del algoritmo, propuesto por Rinc on y col. (2009).

REFERENCIAS BIBLIOGR FICAS

- Elias, P. (1975). Universal codeword sets and representations of the integers. IEEE Transactions on Information Theory. Volumen 21, edici on 2. (194-203).
- Fraenkel, A. y Klein, S. (1985). Novel compression of Sparse Bit-Strings-Preliminary report. En Apostolico, A. y Galil, Z. (editors). Combinatorial Algorithms on Words. Volumen 12, NATO ASI Series F. Estados Unidos. Springer-Verlag.
- Golomb, S. (1966). Run-length encodings (Corresp.). Information Theory, IEEE Transactions. Volumen 12, n umero 3. (Pp. 399-401).



Huffman, D. (1952). A method for the construction of minimum redundancy codes. Proceedings of the IRE. Volumen 40, edición 9. (Pp. 1098-1101).

Rincón, C.; Rodríguez, D.; Acurero, A.; Bracho, D. y Jakymec, J. (2009). Algoritmo de compresión probabilístico basado en la teoría de la información. Octava Conferencia Iberoamericana en Sistemas, Cibernética e Informática, CISCI 2009. Julio, Estados Unidos.

Rincón, C.; Bracho, D.; Acurero, A. (2012). Efecto del tamaño del alfabeto en el rendimiento de un algoritmo de compresión probabilístico. Revista Enl@ce. Volumen 13, número 3. (Pp. 83-93).

Shannon, C. (1948). A mathematical theory of communication. The Bell Systems Technical Journal. Volumen 27. (Pp. 379-423, 623-656).