

INSTITUTO TECNOLÓGICO DE COSTA RICA

SEGUNDO PROYECTO PROGRAMADO

Ariel Herrera

Saúl Zamora

profesor

M. Sc. Saúl Calderón Ramírez

September 23, 2016

I. INTRODUCCIÓN

En la actualidad, los sistemas de almacenamiento y comunicación digitales requieren de métodos optimizados para el uso de recursos (energéticos, temporales, de espacio, etc). Para satisfacer tal necesidad de manera efectiva, muchas disciplinas han formulado múltiples algoritmos para comprimir y descomprimir la información. La compresión de datos consiste en aplicar algún método que permita reducir el tamaño original de la información. Los algoritmos de compresión sin pérdida son capaces de aplicar una serie de pasos para construir la información comprimida, para luego, cuando la información original necesite ser accesada, se descomprime y recupera la información original completamente idéntica. Los algoritmos de compresión con pérdida en cambio, al implementar la descompresión de la información, no lo gran recuperar el 100% de los datos originales. El algoritmo de Huffman implementado en este proyecto fue propuesto por David A. Huffman en 1952 enfocado en la compresión sin pérdida de datos.

II. ANÁLISIS DEL PROBLEMA

La técnica de Huffman trabaja al crear un árbol binario de nodos, los cuales pueden ser hojas o nodos internos. Al principio, todos empiezan como hojas, las cuales contienen un símbolo, el peso (*frecuencia*) es opcional, y un enlace al nodo padre, lo cual facilita leer el código comenzando de las hojas. Los nodos internos contienen el peso del símbolo, dos enlaces a nodos hijos y un enlace opcional a un nodo padre.

Como una convención, el bit 0 representa el siguiente hijo izquierdo y el bit 1 el siguiente hijo derecho. Un árbol terminado puede crecer hasta tener (n) hojas y ($n - 1$) nodos internos. Un árbol de Huffman que omite los símbolos que no se usan, produce el código con el largo óptimo.

El proceso inicia con las hojas conteniendo los símbolos a representar, luego un nuevo nodo es creado con los nodos con menor probabilidad como hijos, tal que la probabilidad de dicho nodo es igual a la suma de las probabilidades de sus hijos. Con los nodos anteriores mezclados en uno (ya no son considerados), y considerando al nuevo nodo, el proceso se repite hasta que solo quede un nodo: el árbol de Huffman.

III. DISEÑO DE LA SOLUCIÓN

A. El algoritmo de Huffman

El algoritmo de Huffman fue diseñado para comprimir señales digitales. Dichas señales están compuestas por un conjunto finito de símbolos, donde cada uno está definido por una cadena de bits.

El algoritmo de Huffman busca crear para cada símbolo, una cadena de bits o código de Huffman, que al reemplazarse por el *token* correspondiente, reduzca el tamaño total del texto. Para construir el diccionario, el algoritmo analiza todo el texto, para construir un diccionario de frecuencias de aparición, el cual defina una entrada por *token*, donde dicho *token* es la llave y el valor de entrada es definido por la cantidad de veces que el *token* aparece en el texto.

El algoritmo busca generar el código más corto para el símbolo más recurrente. Para ello es necesario asegurarse que

la codificación de todos los tokens no sea ambigua, es decir, que sea posible reconstruir el texto original a partir de la señal codificada. Para ello, se utiliza el árbol binario como estructura de datos.

B. El árbol binario

El árbol binario es una estructura de datos que está compuesta por un conjunto de nodos. Un nodo se entiende como una estructura que puede contener cualquier tipo de dato en su interior. Los datos dentro del nodo forman la *etiqueta* del mismo. Un nodo puede estar enlazado como máximo, a dos nodos (de ahí el nombre, *árbol binario*), los cuales están en un nivel jerárquico inferior.

Un nodo tiene como propiedades su etiqueta, un hijo izquierdo y un hijo derecho. Tales propiedades definen sus operaciones básicas: *crearNodo(nodo)*, que recibe como mínimo la etiqueta como argumento y las inserciones de los nodos izquierdos y derechos *insertarHijoIzquierdo(nodo)* e *insertarHijoDerecho(nodo)* respectivamente.

1) Estructura básica de un árbol binario:

- Raíz: es el único nodo que no descende de ningún otro. Es el nodo jerárquicamente superior.
- Nodos hoja: se definen como aquellos nodos que no tienen descendientes.

C. Huffman y el árbol binario

El algoritmo consiste en construir el árbol binario *de abajo hacia arriba* para definir el código de cada token. Al terminar de construir el árbol, existirá un nodo por cada token, y que también contendrá la cantidad total de repeticiones de todos los nodos hijos, en cada nodo.

Los pasos para construir el árbol es la siguiente:

- 1) Crear un nodo por cada token, incluyendo en su etiqueta el número de repeticiones en el texto. Agregar todos los nodos a una lista.
- 2) Tomar y remover los dos nodos de menor frecuencia de la lista y unirlos en un nuevo nodo, el cual tendrá los dos nodos con menor frecuencia como hijos. Este nuevo nodo tendrá como etiqueta un token nulo, y como cantidad de apariciones, la suma de la cantidad de apariciones de los nodos hijos. El nuevo nodo se inserta en la lista de nodos.
- 3) Repetir el paso 2 hasta que exista un solo nodo en la lista. Cuando exista un único nodo en la lista de nodos, el mismo se toma como raíz del árbol de Huffman.

El árbol de Huffman es utilizado para construir el código de Huffman para cada token. Dada la naturaleza del algoritmo, los tokens con mayor cantidad de apariciones en el texto son incluidos de último en el árbol, por lo que estarán más cerca de la raíz del árbol binario.

IV. PRUEBAS

V. REFERENCIAS

REFERENCES

- [1] Mamta Sharma. Compression using Huffman coding. *IJCSNS International Journal of Computer Science and Network Security*, 10(5):133141, 2010.