

Enhancement in File Compression Using Huffman Approach

Suman

Dept. of CSE, Delhi Institute of technology and Management, DCRUST , Murthal, India

Abstract- In this paper, we are showing how we can enhance the file compression by using Huffman approach. The main aim of the thesis is to utilize the concept of type casting and data normalization to show that it is good practice to this concept along with Huffman Coding. We are enhancing the concept of type casting and data normalization with Huffman Coding so as to increase the compression efficiency of Huffman Coding. Enhancement will be done in such a way that compressed data would be decompressed there would be no loss of data i.e. integrity of Huffman Coding would remain.

Keywords – *Huffman approach, Compression, Decompression, Typecasting and Data normalization.*

I. INTRODUCTION

In computer science and technology, Huffman coding is an entropy encoding algorithm used for lossless data compression. The term refers to the use of a variable-length code table for encoding a source symbol (such as a character in a file) where the variable-length code table has been derived in a particular way based on the estimated probability of occurrence for each possible value of the source symbol. It was developed by David A. Huffman while he was Ph.D. student at MIT, and published in the 1952 paper "A Method for the Construction of Minimum-Redundancy Codes".

A. Huffman Approach

The technique works by creating a binary tree of nodes. These can be stored in a regular array, the size of which depends on the number of symbols, n . A node can be either a leaf node or an internal node. Initially, all nodes are leaf nodes, which contain the symbol itself, the weight (frequency of appearance) of the symbol and optionally, a link to a parent node which makes it easy to read the code (in reverse) starting from a leaf node. Internal nodes contain symbol weight, links to two child nodes and the optional link to a parent node. As a common convention, bit '0' represents following the left child and bit '1' represents following the right child. A finished tree has up to n leaf nodes and $n - 1$ internal nodes. A Huffman tree that omits unused symbols produces the most optimal code lengths. The process essentially begins with the leaf nodes containing the probabilities of the symbol they represent, and then a new node whose children are the 2 nodes with smallest probability is created, such that the new node's probability is equal to the sum of the children's probability. With the previous 2 nodes merged into one node (thus not considering them anymore), and with the new node being now considered, the procedure is repeated until only one node remains, the Huffman tree. Huffman's procedure creates the optimal code for a set of symbols and probabilities' subject to the constraints that the symbols be coded one at a time . After the code has been created coding or Decoding is accomplished in a simple look up table manner. The code itself is an instantaneous uniquely decodable block code. It is called a block code because each source symbol is mapped into a fixed sequence of code symbols.

B. Procedure of Huffman Approach

The Huffman algorithm generates the most efficient binary code tree at given frequency distribution. Prerequisite is a table with all symbols and their frequency. Any symbol represents a leaf node within the tree.

Compression

The following general procedure has to be applied:

- search for the two nodes providing the lowest frequency, which are not yet assigned to a parent node
- couple these nodes together to a new interior node
- add both frequencies and assign this value to the new interior node

The procedure has to be repeated until all nodes are combined together in a root node.

Decompression

For decoding the Huffman tree is passed through with the encoded data step by step. Whenever a node not having a successor is reached, the assigned symbol will be written to the decoded data.

C. Type Casting

Typecasting and coercion refer to different ways of, implicitly or explicitly, changing an entity of one data type into another. This is done to take advantage of certain features of type hierarchies or type representations. One example would be small integers, which can be stored in a compact format and converted to a larger representation when used in arithmetic computations. In object-oriented programming, type conversion allows programs to treat objects of one type as one of their ancestor types to simplify interacting with them.

Benefits of Using Type Casting

Typecast always returns the same number of bytes in the output Y as was in the input X. For example, casting the 16-bit integer 1000 to uint8 with typecast returns the full 16 bits in two 8-bit segments (3 and 232) thus keeping its original value ($3*256 + 232 = 1000$). The cast function, on the other hand, truncates the input value to 255.

D. Database Normalization

Data normalization is a process in which data attributes within a data model are organized to increase the cohesion of entity types. In other words, the goal of data normalization is to reduce and even eliminate data redundancy, an important consideration for application developers because it is incredibly difficult to store objects in a relational database that maintains the same information in several places. Table 5.1 summarizes the three most common forms of normalization (First normal form (1NF), Second normal form (2NF), and Third normal form (3NF)) describing how to put entity types into a series of increasing levels of normalization. Higher levels of data normalization are beyond the scope of this article. With respect to terminology, a data schema is considered to be at the level of normalization of its least normalized entity type. For example, if all of your entity types are at second normal form (2NF) or higher then we say that your data schema is at 2NF.

Table 1. Data Normalization Rules.

Level	Rule
First normal form (1NF)	An entity type is in 1NF when it contains no repeating groups of data.
Second normal form (2NF)	An entity type is in 2NF when it is in 1NF and when all of its non-key attributes are fully dependent on its primary key.
Third normal form (3NF)	An entity type is in 3NF when it is in 2NF and when all of its attributes are directly dependent on the primary key.

II. PROPOSED ALGORITHM

In Matlab we have implemented Huffman coding in a two step procedure the first step is converting normal file to Huffman file and getting the compressed code outside. In next step or second step we convert that compressed code to back to original text code.

A. Proposed Compression Method

The proposed technique can be clearly understood by the block diagram below:

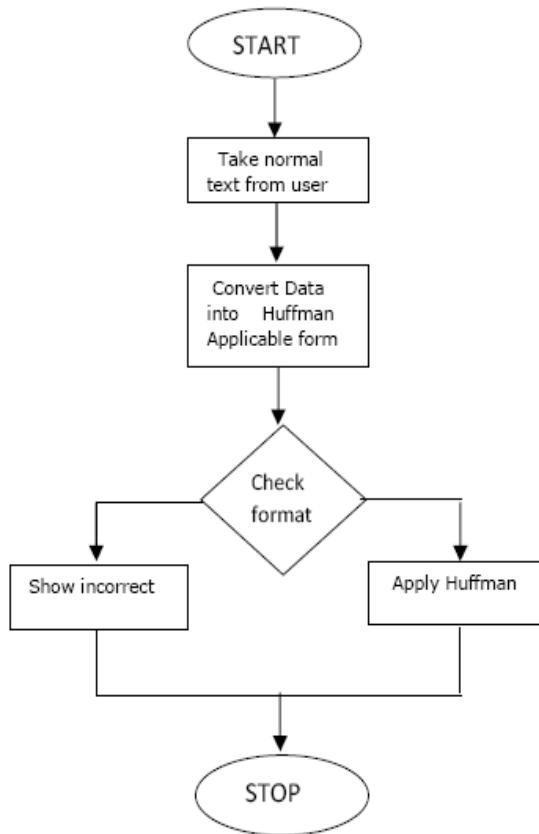


Fig. 1 Proposed Compression Diagram

In the above algorithm we can see that first our Matlab code is taking data from the user. After that it is converting that data to UINT8 format which can be compressed using Matlab if by any chance data is not in UINT8 format then it cannot be compressed. After it is converted and we apply Huffman compression and get a compressed code.

Compression

The following general procedure has to be applied:

- search for the two nodes providing the lowest frequency, which are not yet assigned to a parent node
- couple these nodes together to a new interior node
- add both frequencies and assign this value to the new interior node

The procedure has to be repeated until all nodes are combined together in a root node.

B. Proposed Decompression Method

The proposed technique can be clearly understood by the block diagram below:

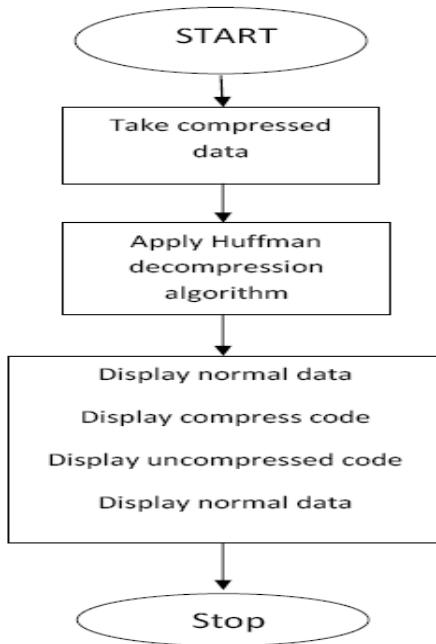


Fig. 2 Proposed Decompression Diagram

Decompression

For decoding the Huffman tree is passed through with the encoded data step by step. Whenever a node not having a successor is reached, the assigned symbol will be written to the decoded data.

III. Result

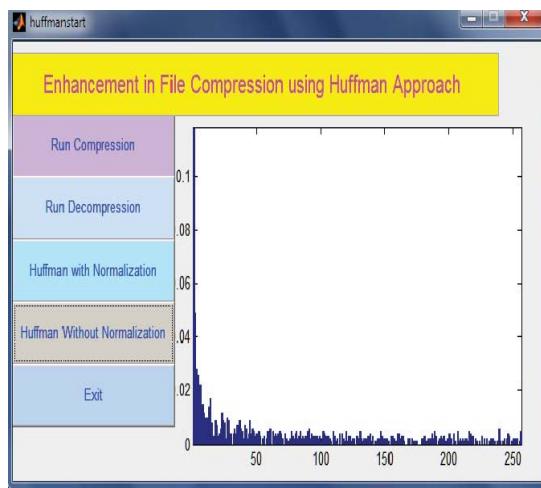


Fig. 1 File compression without normalization.

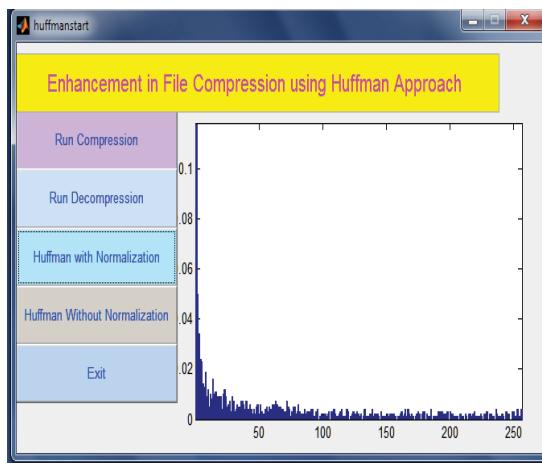


Fig. 2 File compression with normalization.

Source File	Original Size(Kb)	After Compression (kb)	Compression Ratio(kb)
Demo.txt	29	16	0.5517
Demo1.pptx	61	59	0.9672
Demo2.docx	164	163	0.9939

Table. 1 Showing compression ratio

IV. CONCLUSION

With this we can conclude that our algorithm is doing more compression by using normalized data the reason behind this is because it is very easy to distribute frequency in normalized data then normal data due to this reason we are achieving more compression ratio. By using Matlab as development language we are having so advantages which are:

MATLAB is general-purpose mathematical software that is user friendly and very useful for data analysis, simulations, presentation of results, and more. The software runs on both UNIX and WINDOWS.

MATLAB is an interpreter, which means that it performs every command line immediately after it is entered. So you can either enter your commands line by line (prompt mode), or prepare a script file first, and then run a program.

Other advantage of MATLAB is the natural notation used. It looks a lot like the notation that you encounter in a linear algebra course. This makes the use of the program especially easy and it is what makes MATLAB a natural choice for numerical computations.

V. FUTURE SCOPE

Data compression is most consideration thing of the recent world. We have to compress a huge amount of data so as to carry from one place to other or in a storage format. That is why data has to compress. This proposed compression technique has improved the efficiency of file (like .txt, .docx, .pptx) compression using Huffman Approach with the concepts of Typecasting and Data Normalization. For the further research in future you can try this approach to compress file like image etc. to improve the compression efficiency.

ACKNOWLEDGMENT

First of all I am thankful to almighty GOD to help me out in every odds and buts. There are numerous people without whom this dissertation might never have been completed. Among those who deserve credit and our heartfelt gratitude is the most important **Mr. Rakesh Chawla** for his valuable guidance and support. Without his this work would not have been possible. I am thankful for their constant encouragement and guidance. Last but not the least I would like to thank the faculty members of the Department of Computer Science Engineering for being very supportive to me throughout these years. I am also thankful to my friends and classmates for their valuable suggestions.

REFERENCES

- [1] Enhancement in File compression using Huffman approach, IJRASET International Journal For Research In Applied Science And Engineering Technology , Vol. 2 Issue II, Feb. 2014.
- [2] Ternary Tree & A new Huffman Technique, IJCSNS International Journal of Computer Science and Network Security, Vol.10 N0.3, March 2012.
- [3] Typecasting, Legitimation, and Form Emergence: A Formal Theory Greta Hsu Univ. of California at Davis Michael T. Hannan Stanford University László Pólos Durham University Running head: Typecasting, Legitimation, and Form Emergence March, 2012.
- [4] A Study and implementation of the Huffman Algorithm based on Condensed Huffman Table, 2008 International Conference on Computer Science and Software Engineering.
- [5] A Method for the Construction of Minimum-Redundancy Codes David A. Huffman, Associate, IRE, (1952).
- [6] Typecasting, Legitimating: A Formal Theory Greta Hsu Univ. Of California at Davis Michael t. Hannan Stanford University.
- [7] A Method for the Construction of Minimum-Redundancy Codes DAVID A. HUFFMAN, ASSOCIATE, IRE.
- [8] A study and implementation of the Huffman Algorithm based on condensed Huffman table, 2010 International Conference on Computer Science and Software Engineering.
- [9] Typecasting, Legitimation, and Form Emergence: A Formal Theory Greta Hsu Univ. of California at Davis Michael T. Hannan Stanford University László Pólos Durham University Running head: Typecasting, Legitimation, and Form Emergence, March, 2012.
- [10] Introduction to Data Compression, Khalid Sayood, IJCSNS International Journal of Computer Science and Network Security, Vol.10 N0.3, March 2012.
- [11] D.A. Huffman, "A Method for the Construction of Minimum Redundancy Codes", Proceedings of the I.R.E., September 1952, pp 1098–1102. Huffman's original article.
- [12] Database System Concepts By SILBERSCHATZ, KORTH, SUDARSHAN McGraw-Hill Higher Education, ISBN NO. 0-07-120413-X.