LOSSLESS DATA COMPRESSION ANALYSIS OF COMPUTER ALGORITHMS

Nitheesh Edla
Department of computer science
University of North Texas
(UNT) Denton, Texas, USA
nitheeshedla@my.unt.edu

Pragnya Gannapureddy Department of computer science University of North Texas (UNT) Denton, Texas, USA pragnyagannapureddy@my.unt. edu

1. Objective

Loss less data compression is a technique that reduces the size such that a decompression function can restore the original file exactly with no loss of data. The main objective of the project is to decrease the utilization of resources, such as transmission bandwidth or disk space. The steps included in this are encoding the input file and convert it into binary codes and then again decoding it. Using greedy approach, we are implementing this loss less data compression algorithm.

2. Motivation:

We discovered that lossless data compression is extremely useful when compressing data by encoding and decoding, and we also learned about a few applications that use this technique. We are attempting to make the file as small as possible in comparison to the applications, and this project will require a significant amount of learning.

3. Background:

After looking at a variety of compression techniques such as Shannon-Fano Coding, Arithmetic Coding, Huffman Coding, and a few audio compression algorithms, we decided to create an algorithm that is similar to Huffman Coding but uses fewer bytes after compression. We aimed to achieve our desired outcome by

doing the polar opposite of Huffman coding in this approach.

4. Application of Algorithms:

One quite simple approach of compression is run-period encoding, in which massive runs of consecutive same data values are changed with an easy code with the data value and period of the run. This is an instance of lossless records compression. It is regularly used to higher use disk area on workplace computers, or higher use the relationship bandwidth in a pc network. For symbolic data like spreadsheets, text, executable programs, etc., loss lessness is important due to the fact converting even a single bit cannot be tolerated (besides in a few restricted cases).

Generally, for lossless data compression, we use Huffman coding which is a greedy approach technique where each character gets ASCII value and that value is converted into 8-bit binary codes. But the disadvantage of using Huffman coding is it is slower for reading and writing files. Our application has an improvised version of Huffman, making faster reading and writing input files.

5. ALGORITHM:

We'll use the Greedy approach to create an algorithm that performs four jobs, each of which will help us reach our aim of lossless data compression and encoding.

The four tasks are as follows:

- 1. Obtain and save the frequencies of all the distinct characters found in the file.
- 2. Create a tree and store all of the characters according to the frequency created before, assigning binary codes to each leaf and parent node.
- 3. Combine all of the binary values along the way to encode all of the characters in the tree. (For example, if we have a character at the third level and on the right, the binary value of that character would be 01 (if the character's parent is a left node) or 11 (if the character's parent is a right node) (if the parent is also a right node).
- 4. Rearrange the binary values by decoding them into characters.

The below function does the frequency calculation and build a tree:

// Creating a priority queue 'priority' to store the present //nodes from the tree; priority_queue<Node*, vector<Node*>, comp> priority;

```
// Creating a leaf node for each character and
pushing it in to the priority queue created earlier.
for (auto item: fre) {
    priority.push(getNode(item.first, item.second,
    nullptr, nullptr));
    }
```

```
while (priority.size() != 1)
// From the queue, remove the two nodes which
has the highest priority (lowest frequency).
Node *left node = priority.top(); priority.pop();
Node *right node = priority.top();
       priority.pop();
//creating a new node, that is created by adding
up the //frequencies of the earlier popped nodes
and add the new //node to the existing priority
queue.
int sum = left_node->fre + right_node->fre;
priority.push(getNode('\0', sum, left node,
right node));
// root stores a pointer to root of Huffman Tree i.e
top node of the priority queue
       Node* root = priority.top();
// traverse the Tree and store Codes in a map .
Also prints them
unordered map<char, string> lossless;
encode string(root, "", lossless);
cout << "\nGiven string is :" << text << '\n';
       cout << "Binary Codes for the characters
are :\n" << '\n';
       for (auto item: lossless) {
                   cout << item.first << " " <<
               item.second << '\n';
       // printing encoded string
       string str = "";
       for(int i=0;text[i]!='\0';i++)
               str=str+lossless[text[i]];
  cout << "";
   ofstream f;
   f.open("example.txt");
 f \ll str:
```

Lossless data compression and encoding

```
cout << "written to file";
 f.close();
        cout << "\nEncoded string is :\n" << str
<< '\n':
}
       // printing encoded string
        string str = "";
        for(int i=0;text[i]!='\0';i++)
                str=str+lossless[text[i]];
  cout << "";
   ofstream f;
   f.open("example.txt");
 f \ll str:
 cout << "written to file";
 f.close();
        cout << "\nEncoded string is :\n" << str
<< '\n':
}
```

6. Experiments:

6.a Datasets:

We have used several datasets of size 2⁰ to 2¹¹, and the type of the file is text, following is the snapshot of the files we have used in the project.

List of input files:

Name	Date modified	Туре	Size
file1	12/2/2021 12:32 AM	Text Document	1 KB
file2	12/1/2021 3:51 PM	Text Document	2 KB
file3	12/1/2021 3:51 PM	Text Document	4 KB
file4	12/1/2021 3:52 PM	Text Document	8 KB
file5	12/1/2021 3:55 PM	Text Document	16 KB
illie6			32 KB
file7	12/1/2021 3:55 PM	Text Document	63 KB
file8	12/1/2021 3:55 PM	Text Document	125 KB
file9	12/1/2021 3:56 PM	Text Document	250 KB
file10	12/1/2021 4:54 PM	Text Document	489 KB
file11	12/1/2021 4:54 PM	Text Document	977 KB

File data:

File Edit Format View Help
HYKsVCNnYPzi5QS63NhtowYQM+Y5JdZKznnw55Ku3w2eKPMc5EIPpA92UVEm4rrb0FJGPyvygN2
3L8fuUpl.scTib2D4/n2gPvOøVQbe0a/LlpsZIFQEr/pFKK4xJBTHJ9zph9IzzQ3qkR/aBhw8Dte
mjcuzI0nE7SdPbxOrJhBCFiXV51eXNFnL1Cr4M2vB703KwaI0HT8QkL5JtVX4s1W2XBidisidW7
HXENI5DsioB7yx6Jz9MhgruIK8wGujeAuVpWi8gd3qmhrayqhiifc4JRMkLXahkx3052Vm81k60
UPx+tMyHpvfQAT9Ku8zinQLVDzuTx6su+JBox5nZ6BWXrsoptB8dR8g8Bhr95t5onNjIbPqjS17V
ncecnuFchz/LKN2ldAOdTp9ecrIJyQ/8XfxJhsibkodXpfoFJtsZ/LLdeNlCobo95k1MUjbSqoQ
6EyaNeSeYj3ePcN6S4IAiL/NlFoZjFNdU017Ab3x7kUC1LvbZpVRat2hlvEX0L1HAV9TSw6BLRt
k+CCrVWf2KrBpG/O7mlHDTdHrCQYG83WSSWZafH+T3Jxn4WxthPUFWktW0jGwTeCbBY739ZWZH
/RATejwRsRUfCUU93J/EhpAMW02Je7mdiuORtO/xnEjoeFs1pZJ4AXW4G0jbSCEmdMxV81tVlw
Z+0WFA0xHesbE7HYy+yefM6dhWGU9dbz/kStLm8mHEBXXD+45dRUkp/Bv0upBcCeyaAV7V0MIr/
DN/D04gwy0IJHbVQ20/8f7CVGewi9msHQaVC8sk4dNroWNQ05gmBrxJfat4dcMs71h6Fj64ii
7novTpoc1NnWn1dhwt+nicSxhm1k6/1lsncSrcl17a01MncFfaSacsw8ackvs1oTsabhwf/rvbd

6.b Purpose of each experiment:

The purpose of each and every experiment is to get the desired result with less time consuming and we can see that as the size of the file is increasing the time taken for executing is also increasing the size of the output file is also, we can see that the new file generated is encoded into binary and can be seen a slight decrease in size but sometimes based on the number of characters if may differ and for large files the time taken for the execution is so long.

List of output files:

Name	Date modified	Туре	Size
output1	12/8/2021 6:59 PM	Text Document	1 KB
output2	12/8/2021 6:59 PM	Text Document	2 KB
output3	12/8/2021 6:59 PM	Text Document	4 KB
output4	12/8/2021 6:59 PM	Text Document	7 KB
output5	12/8/2021 6:59 PM	Text Document	14 KB
output6	12/8/2021 6:59 PM	Text Document	28 KB
output7	12/8/2021 6:59 PM	Text Document	56 KB
output8	12/8/2021 6:59 PM	Text Document	113 KB
output9	12/8/2021 6:59 PM	Text Document	215 KB
output10	12/8/2021 6:59 PM	Text Document	406 KB
output11	12/8/2021 6:59 PM	Text Document	875 KB

Characters are encoded as below

```
Given string is :data compression algorithm.
Binary Codes for the characters are :

l 11111
h 11110
t 11110
p 11010
i 1100
n 10111
m 1000
o101
e 0100
o 001
d 10110
r 000
g 10011
a 011
```

Output file data:

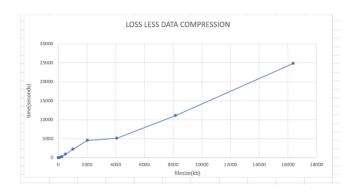
a output5 - Notepad	
File Edit Format View Help	
01111111010110000000011001001011111101111	0011
100000000001010011001001110110110110110	1011
1010100000111000110111010110100101001111	1111
0110001111010011000111111111101111110101	
011000101100001110110000100110100001111010	
100011111101001111001010101011111000000	
11001001101000110011000000100001110011100101	
111010110111010011010000110110011001001	
00111111000110110011101011110000000101111	
001011111010011101000010010111100110100101	
1000010000010001010000010110010100111011010	
000100100101101000010100110110011111100011101111	
00101001011000111010100000100110000000101	
10100101011011010000110000011111111000111010	0111
101011110110110100000010100011101000000	1111

6.c Experimental results:

The complexity of the project is Big O of (log n) 2

<u>The execution time for the project in different experiments is as follows:</u>

File size(kb)	Time			
1	0.006			
2	0.0015			
4	0.039			
8	0.174			
16	0.707			
32	3.112			
64	13.672			
128	72.41			
256	241.6			
500	941.2			
1024	2207.9			
2048	4521.2			
4096	5142			
8192	11087.21			
16384	24856.03			



As you can see from the above table and line graph we can analyze that as the file is increasing in a linear the time is also increasing in the same linear fashion

7. Programming Language:

We used C++ programming language to code the Algorithm.

8. Conclusion:

We have succesfully accomplished the project objective to encode any given text and minimize it's size. Although our project has minor limitations that sometimes may provide wrong outputs, say, the file size may not minimize depending on the number of characters provided. We currently developed this to take text inputs and are working towards accepting inputs as a image or a video.

9. References:

- 1. https://en.wikipedia.org/wiki/Lossless_compression
- 2. https://en.cppreference.com/w/cpp/chrono/hightps://e
- 3. The Data Compression Book 2nd edition By Mark Nelson and Jean-loup Gailly

10. Appendix:

https://github.com/Nitheesh3009/CSCE5150_project