

Lab 3 Report

Compression/Huffman Trees

by Connor McLeod, Matthew Jesus

Mathematical Experiments in Computer Science

November 8, 2023

Section 1:

Data compression is a crucial field where mathematical principles and computer science techniques converge to achieve efficient storage and transmission of information. In this lab, we will explore the close relationship between mathematics and computer science, specifically in the realm of lossless data compression.

Lossless data compression is a technique used to reduce the size of files or data collections without any loss of information. This process involves encoding data using algorithms like Huffman Encoding, and Ziv-Lempel Compression. Both of these algorithms rely heavily on mathematical concepts and data structures, making mathematics a fundamental component of the lab.

One of the primary mathematical aspects in data compression is the concept of frequency analysis. When encoding data using Huffman Encoding, you must first build a frequency chart for all characters in your text. This frequency chart is essentially a mathematical representation of the occurrence of each character in the text. The more frequently a character appears, the shorter its encoded representation in the Huffman tree, resulting in efficient compression. Furthermore, the encoding process itself relies on mathematical concepts of binary representation. The codewords generated during the Huffman encoding are binary sequences, which directly correspond to mathematical concepts of binary numbers. The codewords are assigned to each character in the text, and the efficiency of the compression depends on how these codes are constructed and used.

In "Real Life Compression," mathematics plays a significant role in converting bytes into bits for more efficient storage. This process involves converting eight '0' or '1' characters (bytes) into a single byte using binary arithmetic. The mathematical operations involved in this conversion ensure that the compression is achieved optimally.

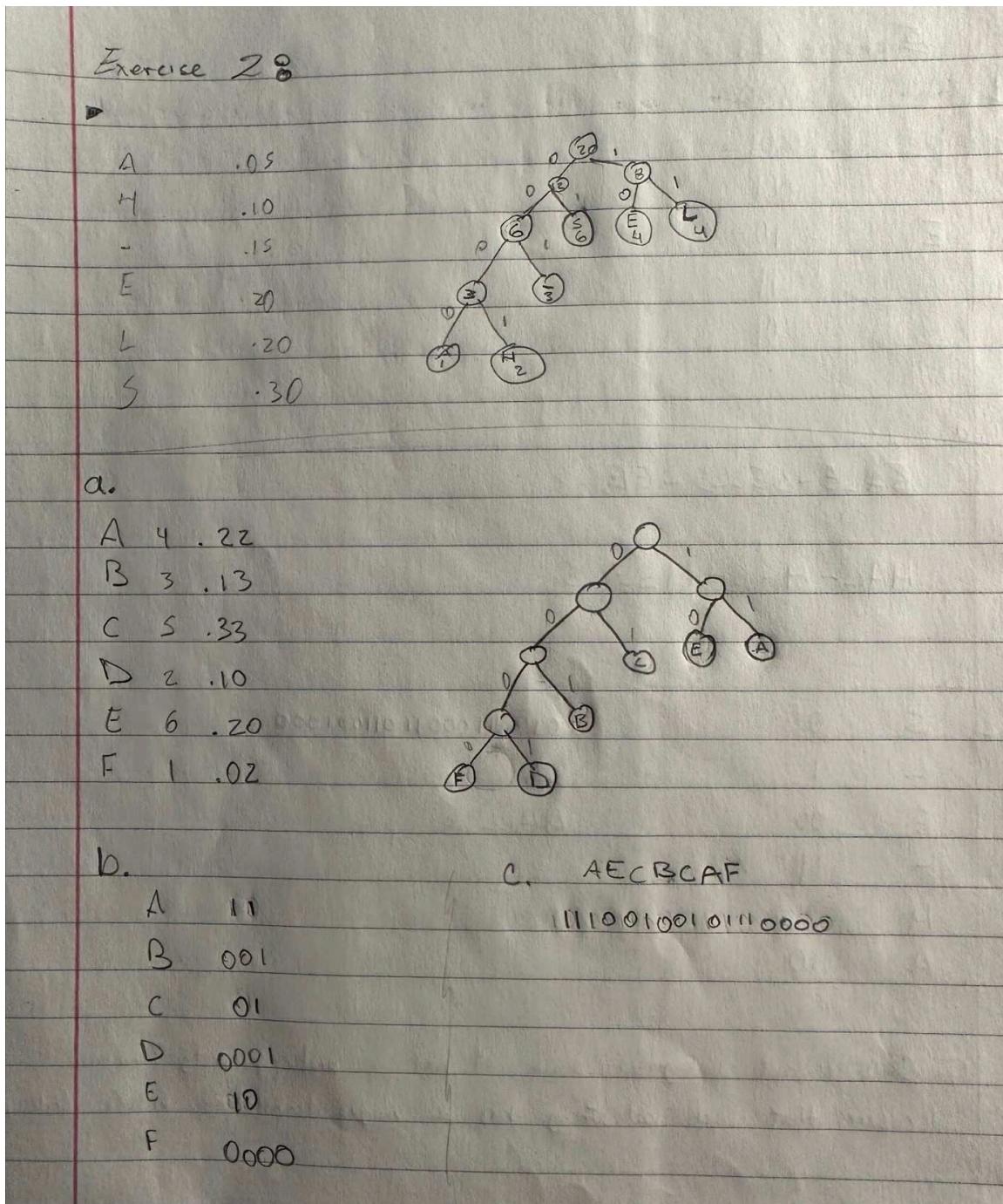
Computer science, on the other hand, provides the tools and techniques to implement these mathematical concepts effectively. Students are required to write programs to encode and decode text using the Huffman algorithm. This process involves designing data structures and algorithms to manipulate data efficiently, which is at the core of computer science.

Section 2:

Exercise 1:

Exercise 1		
Code 1	A	0000
	H	0001
	-	001
	E	10
	L	11
	S	01
a.		
SALE-SELLS-EEELS		
HAL-HAS-HEELS		
b.		
Code 2	S	0
	L	1
	E	00
	-	11
	H	01
	A	10
C. Code 1 isn't ambiguous since it has a prefix-property, which ensures that each bit string has a unique translation unlike Code		

Exercise 2:



Program 1:

The purpose of program one, is to take in a normal text file and output a 'compressed' version of the file that they just provided. The logic behind this process is that the program will create a table that contains the number of times that each letter or symbol was used in the file that was provided. Using that frequency table, the program builds something that is called a

Huffman tree. Using a Huffman tree allows the program to assign a code to each possible letter without any repetition. Once the Huffman tree has been built, the program creates an array of codewords for each possible character and then uses the codewords to encode the originally provided file. However, this is a program that lays down the foundation work to be able to compress files and this program actually makes the files larger because each character which is represented by one byte, actually turns into several different characters which each take up 1 byte of storage.

The photos showcase the file before compression (left), the codewords table (middle), and finally the compressed file (right)

1 THE TRAGEDY OF HAMLET, PRINCE OF DENMARK	1 10
2	2 0100111100
3	3 110010111100000
4 by William Shakespeare	4 11001010111100001
5	5 00110111011
6	6 0011001111111
7	7 0011001111110
8 Dramatis Personae	8 011110
9	9 001100110
10 Claudius, King of Denmark.	10 1100000
11 Marcellus, Officer.	11 11001010111110001
12 Hamlet, son to the former, and nephew to the present king.	12 11001011111101
13 Polonius, Lord Chamberlain.	13 001100101
14 Horatio, friend to Hamlet.	14 110010110
15 Laertes, son to Polonius.	15 001001100
16 Voltimand, courtier.	16 001100100
17 Cornelius, courtier.	17 0100110100
18 Rosencrantz, courtier.	18 0011001011
19 Guildenstern, courtier.	19 000100001
20 Osric, courtier.	20 1100011000
21 A Gentleman, courtier.	21 000100101
22 A Priest.	22 2100101010
23 Marcellus, officer.	23 2200101010
24 Bernardo, officer.	24 230010111
25 Francisco, a soldier	25 24001011111110001
26 Reynaldo, servant to Polonius.	26 250010110
27 Players.	27 2600100100
28 Two Clowns, grave-diggers.	28 2700100110
29 Fortinbras, Prince of Norway.	29 280010110111101
30 A Norwegian Captain.	30 2900100000
31 English Ambassadors.	31 30001011101
32	32 0011001110

Program 2 & Exercise 3:

For program 2, the objective was to be able to take in an encoded file (as created in program 1) and decoding it so that you can read it as it was before. The program works by taking the frequency table that was made in the first program and turning it back into a Huffman tree, because as long as its logic is the same as when it was encoded, you will end up with the

exact Huffman tree. Once the Huffman tree has been built, the code will then use the Huffman tree in order to make a codeword table again which finally is used to decode the encoded file.

Below you will see a encoded file on the left that has been encoded by program 1, and on the right you will see the same file but decoded by program 2. The files shown will be Hamlet, The Constitution, and finally the Declaration of Independence.

THE TRAGEDY OF HAMLET, PRINCE OF DENMARK

by William Shakespeare

Dramatis Personae

Claudius, King of Denmark.
Marcellus, Officer.
Hamlet, son to the former, and nephew to the present king.
Polonius, Lord Chamberlain.
Horatio, friend to Hamlet.
Laertes, son to Polonius.
Voltemand, courtier.
Cornelius, courtier.
Rosencrantz, courtier.
Guildenstern, courtier.
Osric, courtier.
A Gentleman, courtier.
A Priest.
Marcellus, officer.
Bernardo, officer.
Francisco, a soldier
Reynaldo, servant to Polonius.
Players.
Two Clowns, gravediggers.
Fortinbras, Prince of Norway.
A Norwegian Captain.
English Ambassadors.

THE UNITED STATES CONSTITUTION

(See Note 1)

We the People of the United States, in Order to form a more perfect Union, establish Justice, insure domestic Tranquility, provide for the common defence, promote the general Welfare, and secure the Blessings of Liberty to ourselves and our

Posterity, c

Article. I.

Section 1.
All legislative Powers herein granted shall be vested in a Congress of the United States, which shall consist of a Senate and House of Representatives.
Section 2.
Clause 1: The House of Representatives shall be composed of Members chosen every second Year by the People of the

Clauses 1 & 2: No Person shall be a Representative who shall not have attained to the Age of twenty five Years, and been seven years a Member of the most numerous Branch of the State Legislature.

Clause 3: Representatives and direct Taxes shall be apportioned among the several States which may be included within this Union, according to their respective Numbers, which shall be determined by adding to the whole Number of free Persons, including those bound to Service for a Term of years, and excluding Indians not taxed, three fifths of all other Persons. (See Article I, Section 2) The actual Enumeration shall be made within ten Years after the first Meeting of the Congress of the United States, and within every subsequent Term of ten Years, in such Manner as they shall by Law direct. The Number of Representatives

Program 3:

For program 3, the objective was to take the output of program 1 and actually compress it by taking all of the ones and zeroes out of the given file and then taking eight at a time and assigning them to one byte. This will decrease the size of the file because previously, each one and zero was automatically given a size of one byte so by turning eight of them into one, we see a drastic decrease in size. The final step of the program is it finally writes all of the bytes into a file for the user's use. Something that did however surprise me during this lab is that Hamlet and the Constitution compressed down to the same size despite starting with differing sizes. I did also double check this and it ran the same both times.

	Text file	Winzip	Huffman
Hamlet	191,726 bytes	69,800 bytes	20,277 bytes
Constitution	36,267 bytes	11,100 bytes	20,277 bytes
Declaration of Independence	8,487 bytes	3,380 bytes	4,548 bytes

--	--	--	--	--

--	--	--

Mathematical Questions:

Exercise 4:

Huffman's Algorithm produces a code with the prefix property by constructing a binary tree in a way that shorter codes are assigned to shorter symbols. This makes sure that no code assigned to a symbol is a prefix of another symbol's code.

Exercise 5:

Prove: An optimal prefix code for a file is always represented by a full binary tree...*

Let's consider the properties of an optimal prefix code:

1. Shorter codes are assigned to more frequent symbols.
2. Longer codes are assigned to less frequent symbols.

In a non-full tree, we have a non-leaf node "N" with a single child. This child must be a leaf node because all other non-leaf nodes have two children in a binary tree. Let's call this leaf node "L." "L" represents a symbol with a lower frequency because, by definition of an optimal prefix code, the more frequent symbols receive shorter codes.

Now, we can create a better code by merging node "N" and its child "L" into a new node, which has two children, like a full binary tree. The code for this new node would be the concatenation of the codes for "N" and "L."

This merging process reduces the height of the tree because we remove "N" and "L" from the tree and replace them with a single node. This reduction in height implies that the average code length for the entire tree is reduced because the code length for "N" and "L" is replaced by the code length of the merged node, which is shorter. This contradicts the optimality of the original code because we have found a way to achieve a shorter average code length, which should not be possible for an optimal code.

Therefore, we have reached a contradiction, demonstrating that an optimal prefix code for a file cannot be represented by a non-full binary tree. It must be represented by a full binary tree, where every non-leaf node has two children, to ensure that the code is as efficient as possible.

Exercise 6:

Exercise 6:

Show by means of a counter that a full binary tree may represent a prefix code which is not optimal, even if all the letters with lower frequencies have longer codewords.

Code 1: SHE-SELLS-SEA-SHELLS

Char	Codeword	Frequency	Rel. Freq
A	0000	1	.05
H	0001	2	.10
-	001	3	.15
E	10	4	.20
L	11	4	.20
S	01	6	.30

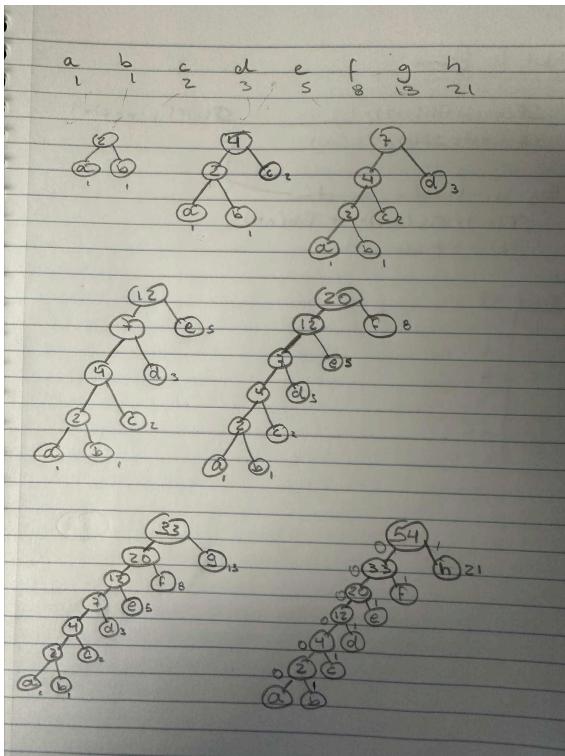
New Codeword

A	11110
H	10
-	11111
E	110
L	1110
S	0

SHE - SELLS - SEA - SHELLS

010110110 01101101100 011011110 01011011011100

Exercise 9:



a.

n	$Fib(n)$	$\text{SumFib}(n) = Fib(1) + Fib(2) + \dots + Fib(n)$
1	1	1
2	1	2
3	2	4
4	3	7
5	5	12
6	8	20
7	13	33
8	21	54
9	34	88

Conjecture:
 $\text{SumFib}(n) = Fib(n+2) - 1$

Base case: $n=1$
 $\text{SumFib}(1) = Fib(3) - 1 = 2 - 1 = 1$

$k = m$
 $\text{SumFib}(m) = Fib(m+2) - 1$

$k = m+1$
 $\text{SumFib}(m+1) = \text{SumFib}(m) + Fib(m+1)$
 $= (Fib(m+2) - 1) + Fib(m+1)$
 $Fib(m+2) = Fib(m+1) + Fib(m) = Fib(m+2) - 1$

Therefore, $\text{SumFib}(n) = Fib(n+2) - 1$ holds for all positive ints
 n

b.

n	$Fib(n)$	$C(n)$
1	1	$0 \cdot 1 = 0$
2	1	$1 \cdot 1 + 1 \cdot 1 = 2$
3	2	$2 \cdot 1 + 1 \cdot 1 = 6$
4	3	$3 \cdot 1 + 1 \cdot 1 = 13$
5	5	$5 \cdot 1 + 1 \cdot 1 = 30$
6	8	$8 \cdot 1 + 1 \cdot 1 = 68$
7	13	$13 \cdot 1 + 1 \cdot 1 = 140$
8	21	$21 \cdot 1 + 1 \cdot 1 = 284$
9	34	$34 \cdot 1 + 1 \cdot 1 = 576$

Prove $C(n) = C(n-1) + \text{SumFib}(n)$, where $C(1) = 0$

For $n=1$
 $C(1) = 0$

For each n , the cost is obtained by adding the cost of the previous tree ($C(n-1)$) with the sum of the first n Fib numbers ($\text{SumFib}(n)$).

$C(n) = C(n-1) + \text{SumFib}(n)$

Substitute $\text{SumFib}(n)$ with $Fib(n+2) - 1$

$$C(n) = C(n-1) + (Fib(n+2) - 1)$$

$$C(n) = (C(n-2) + \text{SumFib}(n-1)) + (Fib(n+2) - 1)$$

Simplify:

$$C(n) = C(n-2) + SumFib(n-1) + Fib(n+2) - 1$$

$$C(n) = C(n-3) + SumFib(n-2) + Fib(n+2) - 1 + Fib(n+1) - 1$$

$$\vdots$$

$$C(n) = 0 + SumFib(2) + Fib(n+2) - 1 + Fib(n+1) - 1 + Fib(n) - n$$

$$C(n) = Fib(4) - 1 + Fib(n+2) + Fib(n+1) + Fib(n) - n$$

(Use $Fib(n+2) = Fib(n+1) + Fib(n)$)

$$C(n) = Fib(n+2) + Fib(n) - 1 - n$$

$$C(n) = Fib(n+2) + Fib(n) - 1 - n$$

$$C(n) = Fib(n+2)$$

$$C(n) = Fib(n) + Fib(n) - 1 - n$$

$$C(n) = Fib(n) + Fib(n) - 1 - n$$

Movie Quiz:

1. What was Fermat's "day job," i.e., how did he make a living? **Fermat's day job was working as a magistrate in the Parlement of Toulouse.**
2. Where did Andrew Wiles live when he worked on Fermat's Last Theorem? **He lived in Princeton University.**
3. The statement of Fermat's Last Theorem was written in the margin of Fermat's copy of Diophantus' book on Arithmetic. According to his notes, why didn't Fermat include his proof? **He didn't include his proof because it was too small to be contained.**
4. How many years elapsed (to the nearest hundred) between the statement of Fermat's Last Theorem and its proof? **About 350 years elapsed between both.**
5. What two mathematical subjects that "live on different planets" are bridged by the Taniyama-Shimura Conjecture? **Elliptic curves and modular forms.**
6. Goro Shimura died in 2019 at the age of 89. What happened to Taniyama? **Taniyama committed suicide in 1958.**
7. The song "One Way or Another" that plays in the movie is by what female 1980's punk vocal artist? (singer or group name is acceptable) **The group's name is Blondie.**
8. What did Andrew Wiles' wife say she wanted for her birthday? **She wanted him to prove the theorem.**
9. For how many years (within 3) did Wiles work on Fermat's Last Theorem? **Wiles worked on Fermat's Last Theorem for 3 years**
10. Which mathematician was not mentioned in the movie The Proof? John Coates*, Andrew Wiles*, Peter Sarnak*, Ken Ribet, Noam Elkies, Sophie Germain, John Conway, Pierre de Fermat*, Goro Shimura*, Barry Mazur, Nick Katz*. **Noam Elkies was not mentioned in "The Proof".**