



Faculty Of Engineering and Technology
Electrical and Computer Engineering Department
Information and Coding Theory ENEE5304
Project Report

Sondos Shahin 1200166

Instructor: Dr. Wael Hashlamoun

Sec. 1

Abstract

The aim of this project is first to generate a variable length random sequence from the symbols (a, b, c, d), given that each symbol has a different probability. In addition to calculating the entropy of the source.

Then, a Lempel-Ziv encoding program was implemented to parse data, and give a number for each resulting phrase. Both the number of bits needed to encode the sequence and the number of bits per symbol were found, in addition to the compression ratio relative to the ASCII code.

Moreover, a Huffman encoding program was implemented to find the codeword for each symbol, and to find the average number of bits per symbol.

Finally, a comparison between the Lempel-Ziv and the Huffman encoding methods is done.

Table of contents

Abstract.....	II
Table of contents.....	III
Table of figures	IV
List of tables.....	IV
Theory	1
Results and Analysis	2
Lempel-Ziv.....	2
Huffman	3
Comparison	4
Conclusion	5
References.....	6
Appendix.....	7
Lempel-Ziv.....	7
Huffman.....	9

Table of figures

Figure 1- Lempel-Ziv results	2
Figure 2- Lempel-Ziv compression ratio	3
Figure 3- Lempel-Ziv avg compression ratio	3
Figure 4- Huffman encoding	4

List of tables

Table 1- comparison table	4
---------------------------------	---

Theory

Lempel-Ziv encoding algorithm is considered one of the fastest compression algorithms. It does not require a priori knowledge of the input data, unlike other algorithms where certain characteristics of data need to be known before compression [1].

Lempel-Ziv compression works by reading a sequence of symbols, then grouping the symbols into different length phrases, and converting these strings into fixed length codes. Since the codes take up less space than the strings they replace, compression will be achieved [2].

Each phrase consists of a previously occurring phrase, followed by the new source output.

On the other hand, Huffman encoding algorithm assigns variable-length codes to input characters. Lengths of the assigned codes are based on the frequencies of corresponding characters, so that symbols with higher probabilities have shorter codewords [3].

In order to use Huffman encoding technique, the probabilities of the source symbols need to be known before starting the encoding.

Results and Analysis

Lempel-Ziv

In figure 1, the results show a 30-symbol randomly generated sequence, that was encoded using the Lempel-Ziv algorithm. The sequence was parsed into phrases, and each phrase was given a number. The resultant codewords are shown in the figure, in addition to the number of bits needed to encode a single symbol, which includes the number of bits used to represent the previously occurring phrases in the dictionary plus the 8-bits used to represent the new ASCII character at the end of the phrase. Finally, the number of bits required to encode the whole sequence, which depends on the number of phrases in the dictionary, and the encoded sequence itself were shown in the figure.

```
sequence is : aababacdaaddbaaabdabababbabbaa
source entropy = 1.6854752972273346
1 a
2 ab
3 aba
4 c
5 d
6 aa
7 dd
8 b
9 aaa
10 bd
11 abab
12 abb
13 abba
codewords: {'a': '0000a', 'ab': '0001b', 'aba': '0010a', 'c': '0000c', 'd': '0000d', 'aa': '0001a', 'dd':
'0101d', 'b': '0000b', 'aaa': '0110a', 'bd': '1000d', 'abab': '0011b', 'abb': '0010b', 'abba': '1100a'}
number of bits per symbol = 5.6
number of bits to encode the sequence = 168
Encoded Sequence: 0000a0001b0010a0000c0000d0001a0101d0000b0110a1000d0011b0010b1100a
```

Figure 1- Lempel-Ziv results

Then, different length sequences were generated, and the compression ratio for the number of bits used for encoding each sequence using Lempel-Ziv algorithm relative to the ASCII code was found. Figure 2 illustrates the results.

Nb is the number of bits used to encode the sequence using Lempel-Ziv algorithm.

sequence length	Nb	compression ratio	bits per symbol
800	3088	0.4825	3.86
100	464	0.58	4.64
1000	3536	0.442	3.536
200	736	0.46	3.68
400	1520	0.475	3.8
2000	6834	0.427125	3.417
50	153	0.3825	3.06
20	85	0.53125	4.25

Figure 2- Lempel-Ziv compression ratio

Then, the average values of the number of bits to encode each sequence (N_{avg}), and the number of bits per symbol, in addition to the compression ratio were found by repeating the previous step 5 times and taking the average of the resultant values. Figure 3 illustrates the results.

sequence length	N _{avg}	compression ratio	bits per symbol
800	2873	0.44040625	3.59125
100	306	0.4335	3.06
1000	3485	0.43435	3.485
200	697	0.45049999999999996	3.485
400	1394	0.45049999999999996	3.485
2000	6579	0.4175625	3.2895
50	170	0.425	3.4
20	68	0.425	3.4

Figure 3- Lempel-Ziv avg compression ratio

Huffman

For this part, a Huffman encoding program was implemented to find a codeword for each source symbol. The codewords depend on the probability of each symbol. Since the symbol 'a' has the highest probability, the program assigned the shortest codeword to it, followed by symbol 'b' which has the second highest probability, and hence was assigned the second shortest codeword. Then both 'c' and 'd' have the same probability so both were given the same codeword length.

The average number of bits per codeword was found. Then for a 100-symbol sequence, the number of bits to encode the sequence using Huffman and using ACSII were found. Figure 4 shows the results.

```

Huffman codewords for each character:
a: 0
c: 100
d: 101
b: 11
Average number of bits per codeword: 1.7000
for a 100 symbol sequence:
number of bits for encoding using ASCII = 800
number of bits for encoding using Huffman = 170.00000000000003
compression ratio using Huffman = 0.21250000000000002

```

Figure 4- Huffman encoding

Comparison

After analyzing the results from all the previous sections, the following results were concluded.

Sequence length = 100 symbol	Encoded sequence size	Compression ratio	Bits per symbol
Lempel-Ziv	504	0.622	5.04
Huffman	170	0.21	1.7
ASCII	800	1	8

Table 1- comparison table

Huffman encoding achieves the most reduction in encoded size, requiring only 170 bits, and a compression ratio of 0.21 and an average of 1.7 bits per symbol. Lempel-Ziv encoding, while not as efficient as Huffman coding for this sequence, still encode much better than ASCII encoding. The Lempel-Ziv algorithm reduces the encoded size to 504 bits, with a compression ratio of 0.622 and an average of 5.04 bits per symbol. Finally, ASCII encoding, which assigns a fixed 8 bits per symbol, results in an encoded size of 800 bits for the sequence, leading to a compression ratio of 1.0 and an average of 8 bits per symbol.

Conclusion

By the end of this project, a clear understanding about the encoding algorithms was gained. And the importance and efficiency of using encoding techniques such as Lempel-Ziv and Huffman algorithms instead of using ASCII encoding in data encoding and transmission was understood.

Huffman encoding algorithm performs a good compression in the cases where the probabilities of source symbols are known. Whereas Lempel-Ziv performs highest compression when the encoded sequence length increases, since the repetition of patterns become more likely. Both algorithms perform a better reduction in the size of encoded sequences compared with encoding using ASCII.

References

- [1] <https://www.studysmarter.co.uk/explanations/computer-science/data-representation-in-computer-science/lempel-ziv-welch/> Accessed 9/6/2024 at 14:00.
- [2] <https://www.geeksforgeeks.org/lzw-lempel-ziv-welch-compression-technique/> Accessed 9/6/2024 at 14:00.
- [3] <https://www.geeksforgeeks.org/huffman-coding-greedy-algo-3/> Accessed 11/6/2024 at 19:00.

Appendix

Lempel-Ziv

```
lempel_ziv.py > ...
1  from random import random
2  import numpy as np
3  import math
4
5  def generate_seq(prob_table, length):
6      #prob_table is the input probability dictionary
7      seq_list = ''
8      s = ''
9      entropy=0
10     while len(seq_list) < length:
11         for k, v in prob_table.items():
12             if len(s) == length:
13                 seq_list = seq_list + s
14                 s = ''
15                 break
16             rn = random()
17             if rn <= v:
18                 s += k
19     return seq_list
20
21 def calculate_entropy(probabilities):
22     return -sum( p * math.log2(p) for p in probabilities)
23
24 def parse_data(sequence):
25     k=1
26     temp = ''
27     for i in sequence:
28         if temp in lz_dict.values() or temp == '':
29             temp = temp+i
30             continue
31         lz_dict[k] = temp
32         #print(k, lz_dict[k])
33         k +=1
34         temp = i
35     return lz_dict, k
36
37
38 def find_phrases(num_of_bits):
39     temp = {}
40     for present_value in lz_dict.values():
41         if len(present_value) == 1:
42             temp[present_value] = '0'.rjust(num_of_bits, '0') + present_value[0]
43         else:
44             for past_key, past_value in lz_dict.items():
45                 if present_value[:len(present_value) - 1] == past_value:
46                     temp[present_value] = str(bin(past_key).replace("0b", "").rjust(num_of_bits, '0')) + present_value[len(present_value) - 1]
47     return temp
48
49 def lempel_ziv(sequence):
50     lz_dict, count = parse_data(sequence)
51     num_bits = math.ceil(math.log2(len(lz_dict.items())))
52     codewords = find_phrases(num_bits)
53     return num_bits, codewords, count
54
55
56 def summarize_info(sequence_lengths, prob_table):
57     print("sequence length \t Nb \t compression ratio \t bits per symbol")
58     for length in sequence_lengths:
59         sequence = generate_seq(prob_table, length)
60         lz_dict = {}
61         num_bits, codewords, count = lempel_ziv(sequence)
62         seq_bits = (num_bits+8)*count
63         compression_ratio = seq_bits / (length*8)
64         print(length, " \t\t\t", seq_bits, "\t\t", compression_ratio, " \t\t", seq_bits/length )
65
```

```

66 def summarize_avg(sequence_lengths,prob_table):
67     print("sequence length \t Navg \t compression ratio \t bits per symbol")
68     for length in sequence_lengths:
69         avg_bits=0
70         for i in range (5):
71             sequence = generate_seq(prob_table,length)
72             lz_dict = {}
73             num_bits, codewords, count= lempel_ziv(sequence)
74             seq_bits = (num_bits+8)*count
75             avg_bits += seq_bits
76         avg_bits= avg_bits / 5
77         compression_ratio = avg_bits / (length*8)
78         print(length , " \t\t\t", seq_bits, "\t\t\t ", compression_ratio, " \t\t\t ", seq_bits/length)
79
80
81 prob_table = {'a': 0.5, 'b': 0.3, 'c': 0.1, 'd': 0.1}
82 sequence = generate_seq(prob_table,30)
83 print("sequence is : " ,sequence)
84 entropy = calculate_entropy([0.5,0.3,0.1,0.1])
85 print("source entropy = ",entropy)
86 lz_dict = {}
87 num_bits, codewords, count= lempel_ziv(sequence)
88 sequence_length= (num_bits+8)*count
89 print("codewords: ",codewords)
90 print("number of bits per symbol = ", sequence_length/30)
91 print("number of bits to encode the sequence = ", sequence_length)
92 encoded_sequence = ''.join([str(item) for item in codewords.values()])
93 print("Encoded Sequence: " + encoded_sequence + "\n")
94
95 sequence_lengths = {20,50,100,200,400,800,1000,2000}
96 summarize_info(sequence_lengths,prob_table)
97 summarize_avg(sequence_lengths,prob_table)
98

```

Huffman

```
huffman.py > ...
1 class Node:
2     def __init__(self, char, freq):
3         self.char = char
4         self.freq = freq
5         self.left = None
6         self.right = None
7
8 # Function to build the Huffman Tree
9 def build_huffman_tree(char_prob):
10     # Create a list of nodes
11     nodes = [Node(char, prob) for char, prob in char_prob.items()]
12     # Continue until the list contains only one node
13     while len(nodes) > 1:
14         # Sort the list of nodes by frequency
15         nodes = sorted(nodes, key=lambda x: x.freq)
16         # Get the two nodes with the smallest frequencies
17         left = nodes.pop(0)
18         right = nodes.pop(0)
19         # Create a new internal node with these two nodes as children
20         merged = Node(None, left.freq + right.freq)
21         merged.left = left
22         merged.right = right
23         # Add the merged node back to the list
24         nodes.append(merged)
25     # The remaining node is the root of the Huffman Tree
26     return nodes[0]
27
28 # Function to generate Huffman codes from the Huffman Tree
29 def generate_codes(node, prefix="", codebook={}):
30     if node is not None:
31         # If this is a leaf node, it contains a character
32         if node.char is not None:
33             codebook[node.char] = prefix
34         else:
35             generate_codes(node.left, prefix + "0", codebook) # Traverse the left subtree with the prefix '0'
36             generate_codes(node.right, prefix + "1", codebook) # Traverse the right subtree with the prefix '1'
37     return codebook
38
39 def calculate_avg_bits(huffman_codes, char_prob): # Function to calculate the average number of bits per codeword
40     avg_bits = 0
41     for char in huffman_codes:
42         prob = char_prob[char]
43         code_length = len(huffman_codes[char])
44         avg_bits += prob * code_length
45     return avg_bits
46
47 char_prob = {
48     'a': 0.5,
49     'b': 0.3,
50     'c': 0.1,
51     'd': 0.1,
52 }
53 huffman_tree = build_huffman_tree(char_prob)
54 huffman_codes = generate_codes(huffman_tree)
55 print("Huffman codewords for each character:")
56 for char, code in huffman_codes.items():
57     print(f"{char}: {code}")
58 avg_bits_per_codeword = calculate_avg_bits(huffman_codes, char_prob)
59 print(f"Average number of bits per codeword: {avg_bits_per_codeword:.4f}")
60 print("for a 100 symbol sequence: ")
61 print("number of bits for encoding using ASCII = ", 100*8)
62 print("number of bits for encoding using Huffman = ", 100*avg_bits_per_codeword)
63 print("compression ratio using Huffman = ", (100*avg_bits_per_codeword) / (8 * 100))
64
```