



UNIVERSITY OF BERGEN

INFORMATION THEORY
INF242

First mandatory exercise

Student :

Arthur GAUTIER

September 25, 2023



Contents

1	Lempel-Ziv-Welch	2
1.1	Compression	2
1.1.1	Build the dictionary	2
1.1.2	Implementation	2
1.2	Decompression	3
1.3	Compression rate	3
2	Huffman encoding	4
2.1	Node class	4
2.2	Number of occurrences	4
2.3	Huffman coding	4
2.4	Compression rate	5

1 Lempel-Ziv-Welch

1.1 Compression

1.1.1 Build the dictionary

We start by building a dictionary using the code of the ASCII table. We include the higher and lower case letters of the latin alphabet. We also add the higher and lower letters exclusive to the Norwegian language.

```
1 file = open("folktale.txt", "r", encoding="utf-8")
2 data = file.read()
3
4 # Preparing the dictionary :
5 minuscule = [chr(i+97) for i in range(0,26)]
6 majuscule = [chr(i+65) for i in range(0,26)]
7 norwegian_letters_and_space = [chr(198), chr(216), chr(197), chr(230),
8 chr(248), chr(229), chr(32)]
9 characters = minuscule + majuscule + norwegian_letters_and_space
10 dictionary_size = len(characters)
11 dictionary = {characters[i] : i for i in range(0,dictionary_size)}
```

1.1.2 Implementation

We start by initialize variables :

```
1 str = "" # to build sequences of characters
2 output = [] # store the compressed output
```

Then we iterate on each symbol of the data. We start by create a *new_symbol* by concatenating the current *str* with the current *symbol*. We check if this *new_symbol* is contained in the dictionary. If it is the case, it means that we can potentially add more characters to the characters to it and continue searching for a longer sequences. That's why we update *str* to *new_symbol*.

If *new_symbol* is not in the dictionary, we start by output the code for *str* to the *output*. We then add the *new_symbol* to the dictionary with a new index (*dictionary_size* that we immediately increase to 1 for the next potential new symbol). Finally, we update *str* with the current *symbol*.

```
1 for symbol in data:
2     new_symbol = str + symbol
3     if new_symbol in dictionary:
4         str = new_symbol
5     else:
6         output.append(dictionary[str])
```

```
7         dictionary[new_symbol] = dictionary_size
8         dictionary_size += 1
9         str = symbol
10    if str:
11        output.append(dictionary[str])
12
13    print(output)
```

1.2 Decompression

We start by initializing two empty strings : *decompressed_data* (to stock the compressed data) and *str* (to shape the different values of the dictionary).

We start by iterate on each code in the compressed data. If the code is not already in the dictionary, it means that it is a new one : we add this new code to the dictionary (by concatenating *str* with *str[0]*, we extend the dictionary with a new symbol based on the previous sequence).

We then, add the sequence linked to *code* in *decompressed_data* : we build the original data based on the LZW code.

If the *str* variable is not empty, we add a sequence to the dictionary. The new entry is formed by concatenating *str* with the first character of the sequence linked to *code*.

Then we update *str* with the sequence linked to *code* value : it's going to be used as a prefix for the next *code*.

```
1    decompressed_data = ""
2    str = ""
3    next_code = dictionary_size
4    compressed_data = output
5
6    for code in compressed_data:
7        if code not in dictionary:
8            print(code)
9            print(len(str))
10           dictionary[code] = str + str[0]
11           decompressed_data += dictionary[code]
12           if len(str) != 0:
13               dictionary[next_code] = str + dictionary[code][0]
14               next_code += 1
15           str = dictionary[code]
16
17    print(decompressed_data)
```

1.3 Compression rate

The original file is 8200 bytes. The new file after the LZW compression is 5394 bytes. The compression rate is 66%.

2 Huffman encoding

2.1 Node class

The principle of Huffman coding is based on the creation of a tree structure composed of nodes. We need to create a class for the nodes.

```
1 class Node:
2     def __init__(self, value, letter="", NodeLeft=None, NodeRight=None):
3         """constructor"""
4         self.value = value # to store the attributed value of the node
5         self.letter = letter # to store the attributed letter
6         self.NodeLeft = NodeLeft # if it's not a leaf, left node child
7         self.NodeRight = NodeRight # if it's not a leaf, right node child
8
9     def __lt__(self, Node2):
10        """compare two node values"""
11        return self.value < Node2.value
```

2.2 Number of occurrences

We now need to count the number of each value in the data we have.

```
1 def occurrence(data):
2     """count the number of occurrences for each present letters in data"""
3     # extract the different values:
4     values_to_count = []
5     occ = []
6     for value in data:
7         if value not in values_to_count:
8             values_to_count.append(value)
9     for value in values_to_count:
10        number_occurrences = data.count(value)
11        letter = [value, number_occurrences]
12        occ.append(letter)
13    return occ
```

2.3 Huffman coding

Based on this occurrence table, we can now start to build the Huffman encoding. To do that we will use priority queues using the *heapq* library. A priority queue allows three different operation : add a element to the queue, extract an element with the higher key, test if the queue is empty or not. Priority queue answer to all the need of Huffman encoding implementation.

We begin by placing all the various occurrence values previously calculated into a priority queue represented as nodes. Then, as long as there is more than one node remaining

in the queue, we continue to iteratively combine the two highest nodes into a new parent node in each loop, until only a single root node remains in the queue. We return this root.

```
1 def build_huffman(data):
2     """build the huffman tree"""
3     occ = occurrence(data)
4     priority_file = []
5     for element in occ:
6         n = Node(element[1], element[0])
7         heappush(priority_file, n)
8     while len(priority_file) > 1:
9         n1 = heappop(priority_file)
10        n2 = heappop(priority_file)
11        n3 = Node(n1.value + n2.value)
12        n3.NodeLeft = n2
13        n3.NodeRight = n1
14        heappush(priority_file, n3)
15    # Root :
16    return heappop(priority_file)
```

Now that the tree structure is built, we will traverse the tree and assign values to the leaves. We recursively traverse the tree from the root, assigning the value 0 to the left node and 1 to the right node. Then we create a dictionary to attribute a Huffman code to each of the original code.

```
1 def make_huffman_code(data):
2     """Create the code dictionary"""
3     root_node = build_huffman(data)
4     dic = {}
5     code_huffman(root_node, dic, "")
6     return dic
7
8
9 def code_huffman(n, dictionary, code):
10    """recursive function to create the huffman code"""
11    if is_leaf(n):
12        dictionary[code] = n.letter
13    else:
14        code_huffman(n.NodeLeft, dictionary, code + '0')
15        code_huffman(n.NodeRight, dictionary, code + '1')
```

The other functions then allow us to read the original code and assign, for each element of the code, its attribute referenced in the dictionary. This results in a binary code for the original message.

2.4 Compression rate

The original file is 8200 bytes. The new file after the LZW compression + Huffman encoding is 3279 bytes. The compression rate is 40%.