

# Data Analysis and Algorithm

## Practical 6

Write a program to implement Huffman's code algorithm.

Date.: 11-10-21

Name – Yash Vasudeo Prajapati

Rollno - 022

MSc. Computer Science

6) Write a program to implement Huffman's coding.

Theory :- Huffman Coding is a technique of compressing data to reduce its size without losing any of the details. It was first developed by David Huffman.

Example :-

Consider a string "this is an example"

Char	freq	code	size	
s	2	000	$2 \times 3$	6
i	2	001	$2 \times 3$	6
p	2	010	$2 \times 3$	6
a	2	011	$2 \times 3$	6
m	1	1000	$1 \times 4$	4
x	1	1001	$1 \times 4$	4
l	1	1010	$1 \times 4$	4
f	1	1011	$1 \times 4$	4
t	1	1100	$1 \times 4$	4
n	1	11010	$1 \times 4$	4
e	1	11011	$1 \times 4$	4
"	3	111	$3 \times 3$	9
$12 \times 8 = 96 \text{ bits}$				18
				61 bits

we notice that the string  
without encoding would take

216 bits

And after encoding

$$96 + 18 + 61 = 175 \text{ bits}$$

which is not much but with larger  
strings or strings with more frequency  
of same characters a clear difference  
is seen.

7 steps -  
Let see another example  
'a f s d a s d a s d a'

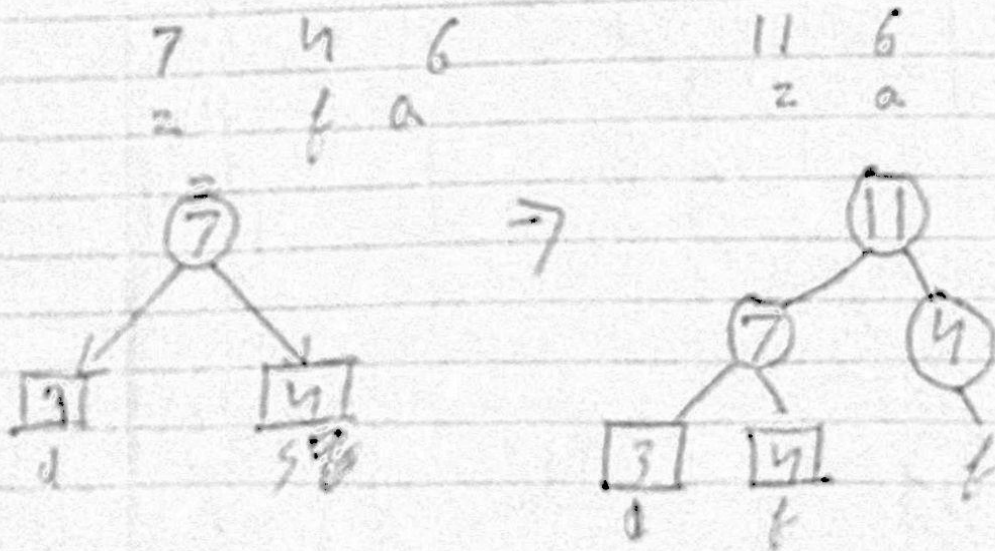
1) Calculate the frequency.

a	f	s	d
↓	↓	↓	↓
6	4	4	3

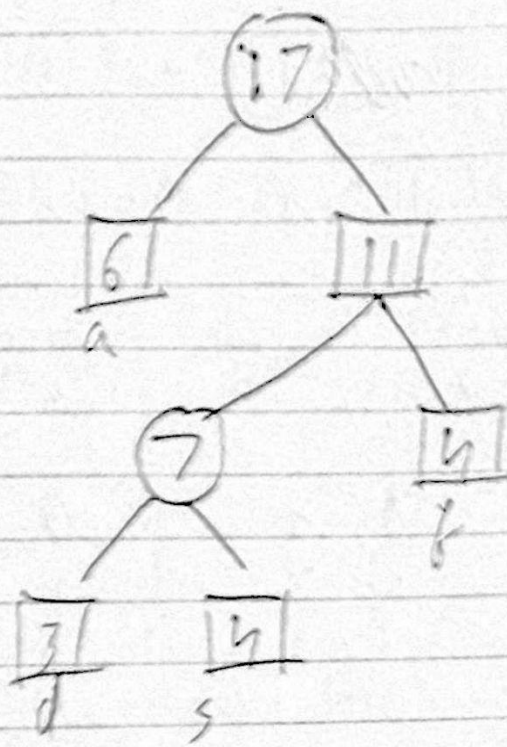
2) Sort in increasing order of frequency

d	s	f	a
3	4	4	6

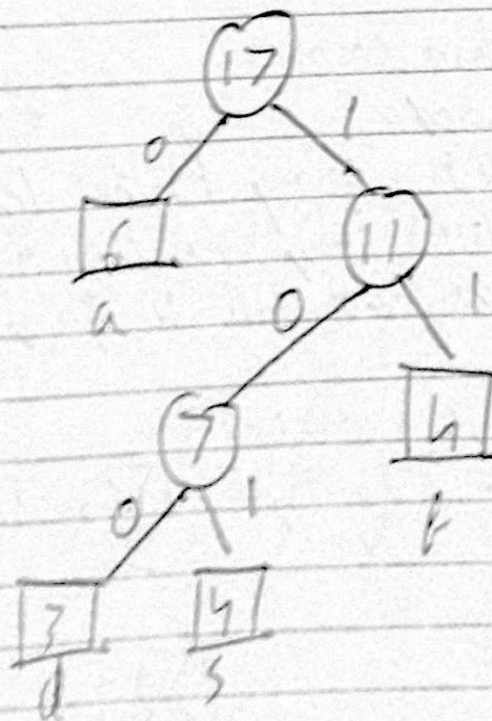
3) Make each character a leaf node, remove the first two nodes & replace with an empty node whose value is the sum of the freq of left & right child & repeat till all nodes are replaced.







4) assign 0 to the left & 1 to right of non-leaf node edge.



We get each code :

a  $\rightarrow$  0  
d  $\rightarrow$  100  
l  $\rightarrow$  11  
s  $\rightarrow$  101

### Algorithm

- 1) Create a list with all unique characters.
- 2) Calculate the frequency of each unique character.
- 3) Repeat for all characters.
  - a) Create a new node.
  - b) Assign the smallest free to the left.
  - c) Assign the smallest free to the right.
  - d) Sum of the above two will be free of new node.
  - e) Insert new node in tree.
- 4) Assign 0 to left & 1 to right nodes recursively.
- 5) Return code.

## Complexity Analysis

To create a min heap tree it takes  $O(n)$

For every node we follow two steps

1. find to 3 binary  $\Rightarrow \log(n) + \log(n)$   
since we are using heap tree.

2. merge  $\Rightarrow \log(n)$

$\Rightarrow$  create a new node & merge 3  
place the new node

thus giving us  $\log(n) + \log(n) + \log(n)$

$$\Rightarrow 3 \log(n) \approx \log(n)$$

The we create tree for all nodes which  
is one less the total node  
 $(n-1) \approx n$

thus we can say the total time taken  
will be  $O(n \log(n))$



## Conclusion

we implement huffman code & find the encoded form of strings & analyse its time complexity.

## Program

```
class NodeTree(object):
    def __init__(self, left=None, right=None):
        self.left = left
        self.right = right

    def nodes(self):
        return (self.left, self.right)

# Main function implementing huffman coding
def huffman_code_tree(node, left=True, binString=''):
    if type(node) is str:
        return {node: binString}
    (l, r) = node.nodes()
    d = dict()
    d.update(huffman_code_tree(l, True, binString + '0'))
    d.update(huffman_code_tree(r, False, binString + '1'))
    return d

# Calculating frequency
string = 'ACDSVDASCASDAW'
freq = {}
for c in string:
    if c in freq:
        freq[c] += 1
    else:
        freq[c] = 1

freq = sorted(freq.items(), key=lambda x: x[1], reverse=True)

nodes = freq

while len(nodes) > 1:
    (key1, c1) = nodes[-1]
    (key2, c2) = nodes[-2]
    nodes = nodes[:-2]
    node = NodeTree(key1, key2)
    nodes.append((node, c1 + c2))
```



```
nodes = sorted(nodes, key=lambda x: x[1], reverse=True)
huffmanCode = huffman_code_tree(nodes[0][0])
print(huffmanCode)
```

Coding.py

```
{'S': '00', 'D': '01', 'W': '1000', 'V': '1001', 'C': '101', 'A': '11'}
```