# HUFFMAN CODES

→ Data can be encoded efficiently using Huffman Codes.

→ It is a technique for compressing data; saving of 20% to 90% depending on the characterstics of the file being compressed.

→ Huffman's greedy algorithm uses a table of the frequencies of occurrence of each character to build up an

→ Example:

Suppose we have $10^5$ characters in a datafile. Normal storage: 8 bits per character (ASCII) means $8 \times 10^5$ bits in the file. But we want to compress the file and store of compactly.

→ Suppose only 6 characters appear in the file:

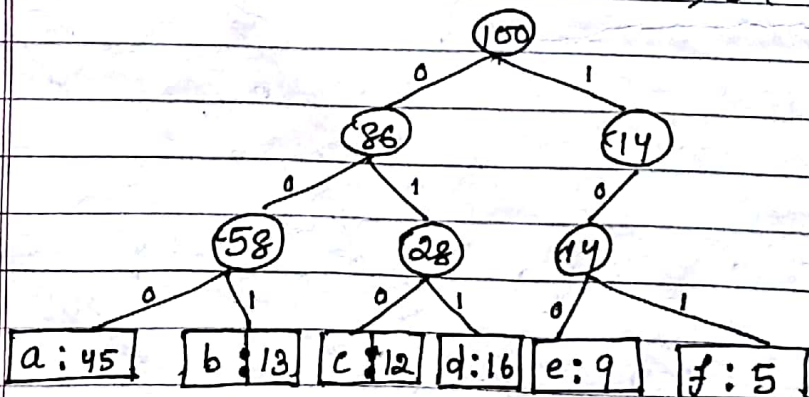|  | a | b | c | d | e | f | Total |
|---|---|---|---|---|---|---|---|
| Frequency | 45 | 13 | 12 | 16 | 9 | 5 | 100 |

How can we represent the data in a compact way.

(i) fixed length code: Each letter represented by an equal number of bits. With a fixed length code at least 3 bits per character:

for example:

| | |
|---|---|
| a | 000 |
| b | 001 |
| c | 010 |
| d | 011 |
| e | 100 |
| f | 101 |

for a file with $10^5$ characters, we need $3 \times 10^5$ bits.



(Not optimal)

## (ii) Variable length code:

Here each character is encoded according to the frequency of the characters.

The more the frequency of characters are short code words and less frequency characters are long code words.

**for Example:**

| | |
|---|---|
| a | 0 |
| b | 101 |
| c | 100 |
| d | 111 |
| e | 1101 |
| f | 1100 |

Number of bits $= (45 \times 1 + 13 \times 3 \times 12 \times 3 + 16 \times 3 + 9 \times 4 + 5 \times 4) \times \dfrac{1}{1000}$

$$= 2.24 \times 10^5 \text{ bits}$$

Thus 224,000 bits required to represent the file.

→ percentage of data compressed is:

$$3,00000 - 224,000 = 76,000$$

$$76,000 / 3,00000 \times 100 = 25\%.$$
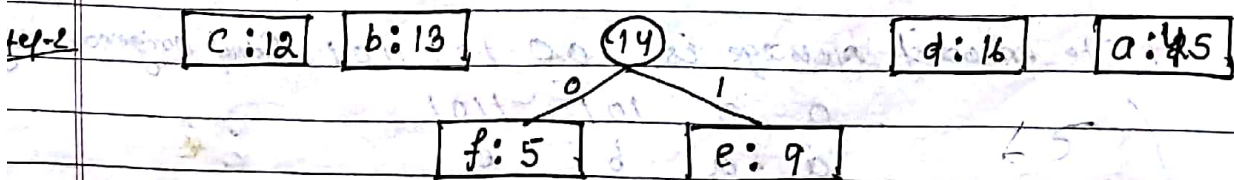
∴ 25% of data is compressed.
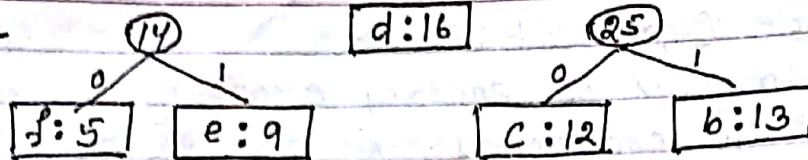
so, this is an optimal character code for this file.

→ **Example:**

step-1

| f : 5 | e : 9 | c : 12 | b : 13 | d : 16 | a : 45 |
|---|---|---|---|---|---|

This algorithm is based on a reduction of a problem with n characters to a problem with n-1 characters. A new character replaces two existing one.
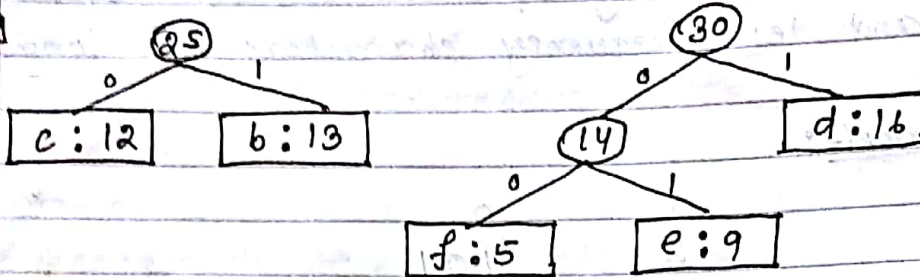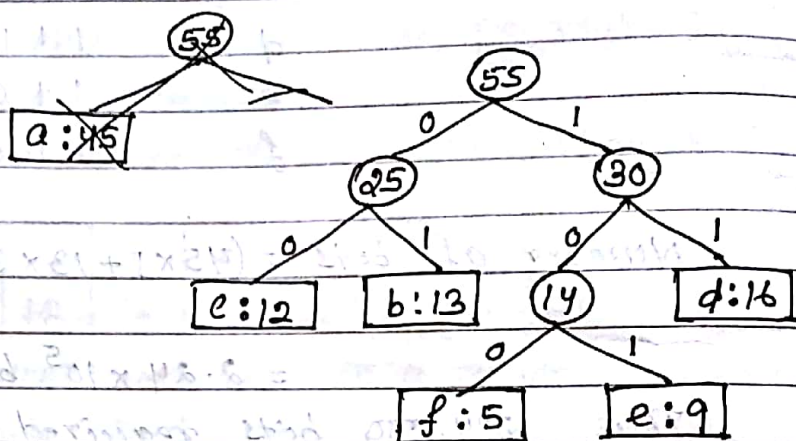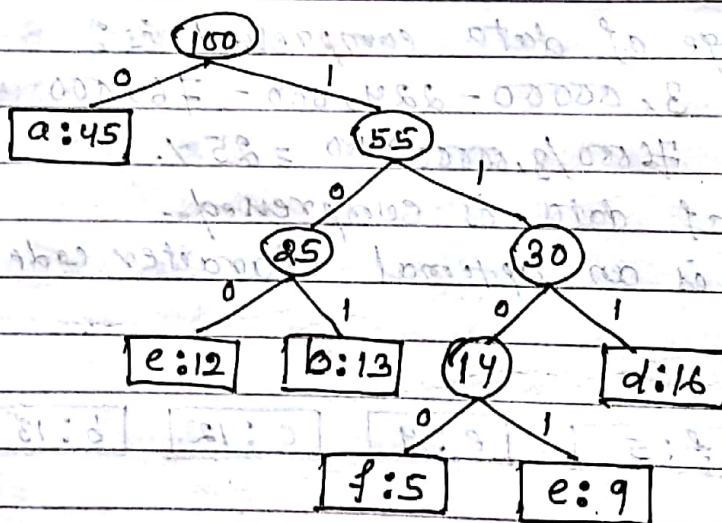
step-2

| c : 12 | b : 13 | (14) | | d : 16 | a : 45 |

| f : 5 | e : 9 |

**step-3**

(14)     d:16     (25)        a: 45

0 / 1          0 / 1

f:5   e:9       c:12   b:13

**step-4**

(25)            (30)        a: 45

0 / 1             0 / 1

c:12   b:13      (14)    d:16

               0 / 1

            f:5   e:9

**step-5**

a: 45        (55)        (55)

     a: 45         0 / 1

               (25)       (30)

         0 / 1      0 / 1

     c:12   b:13    (14)    d:16

                  0 / 1

               f:5   e:9

**step-6**

(100)

0 / 1

a:45     (55)

      0 / 1

   (25)     (30)

  0 / 1    0 / 1

c:12   b:13   (14)    d:16

            0 / 1

          f:5   e:9

| a | b | c | d | e | f |
|---|-----|-----|-----|------|------|
| 0 | 101 | 100 | 111 | 1101 | 1100 |

→ The encoded message is 00 101 1101. find original msg.

      0    0    101    1101

      a    a     b     e

∴ The original message is aabe

66

# Algorithm for Hoffman Code:

HUFFMAN (C)
1. $n \leftarrow |C|$
2. $Q \leftarrow C$                                   $O(n)$
3. for $i \leftarrow 1$ to $n-1$
4.         do allocate a new node $z$
5.            $left[z] \leftarrow x \leftarrow EXTRACT-MIN(Q)$      $(n-1)$
6.            $right[z] \leftarrow y \leftarrow EXTRACT-MIN(Q)$
7.            $f[z] \leftarrow f[x] + f[y]$
8.            $INSERT(Q,z)$
9. return $EXTRACT-MIN(Q)$   ▷ Return the root of the tree

## Complexecity:

→ For complementing Huffman's algorithm, we need a min priority queue.

→ The min priority queue can be implemented by binary min-heap.

→ A set $C$ of $n$ characters are initialized to $Q$ which takes $O(n)$ times. by using the BUILD-MIN-HEAP procedure.

→ Each time it will take 2 nodes whose frequency is minimum and it will delete the 2 nodes from the priority queue and insert a new node whose value is the sum of the frequency of the two deleted nodes which is called merging.

→ The for loop is running $(n-1)$ times

Each heap operation takes $O(\log n)$

So, total time required is $O(n \log n)$