# Encoding

If there is a message having 100 char.

$$a = 5 \quad, \quad b = 9 \quad, \quad c = 12$$

$$d = 13 \quad, \quad e = 16 \quad, \quad f = 45$$

Total characters $\Rightarrow 100$

## ASCII Encoding

Each character takes 8 bits

so

Total bits for this message $= 100 \times 8$

$$= 800 \text{ bits.}$$

# Another method of encoding:-

As there are 6 unique charact

so 3 bits are needed to assign unique code to each character.

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

So Total bits for encoded = $100 \times 3$
message

$= 300$ bits.

* Both above methods are used to assign "Fixed Length codes" to symbols.

In $2^{nd}$ method, message has been compressed (size of bits has been reduced)

800 bits $\longrightarrow$ 300 bits.

* Message can be decoded easily on the basis on scheme used for encoding. because codes are fixed Lengths

# Huffman encoding

It is a way to encode and compress message using "Variable Length codes" to represent characters depending on how frequently they appear.

The ⟨idea⟩ behind this is that
→ characters appeared more frequently should be assigned short codes,
→ while those appear more rarely should be assigned Longer codes

★ Codes assigned to one character is not the prefix of code assigned to any other character.

So it makes sure that there is no ambiguity during decoding.

# Algorithm

assuming index start from 1

HUFFMAN $(N, Symbols[1..N], freq[1..N])$

{

    for $(i \leftarrow 1 \ TO \ N)$

    {

        $t \leftarrow TreeNode \ (symbol[i], freq[i])$

        minHeap. insert $(t, freq[i])$

    }

    for $(i \leftarrow 1 \ TO \ N-1)$

    {

        $x \leftarrow minHeap. remove()$

        $y \leftarrow minHeap. remove()$

        $z \leftarrow new \ TreeNode()$

        $z.left \leftarrow x$

        $z.right \leftarrow y$

        $z.freq \leftarrow x.freq + y.freq$

        minHeap.insert $(z, z.freq)$

    }

    return minHeap.remove()

Time Complexity $= O(n\log n)$

n are unique characters.

---

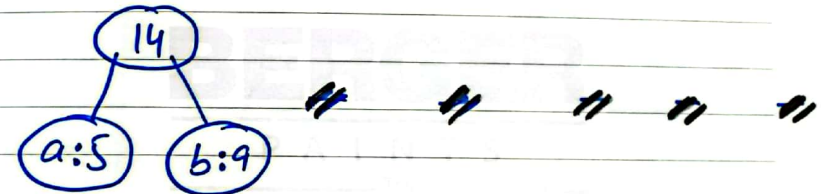Tree finally obtained is the desired Huffman tree.

# Dry Run

Create minHeap or array all nodes of characters in increasing order of their frequencies. (After each step arrange resulting node and Pick minimum 2)
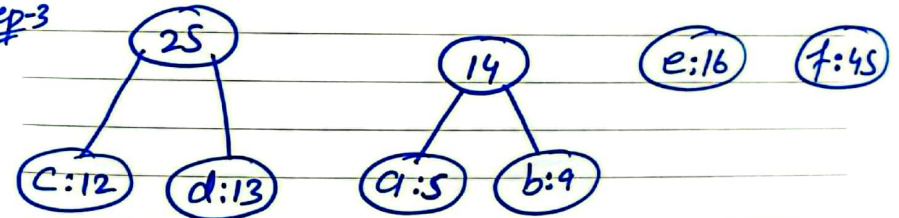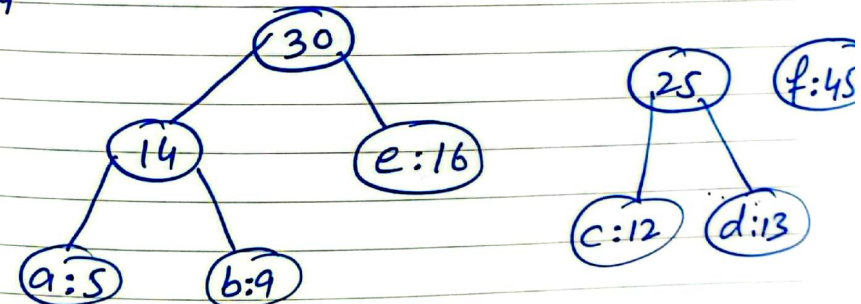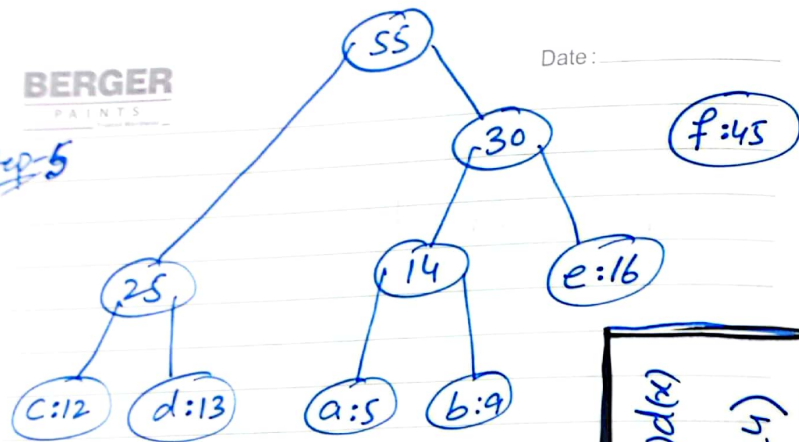
**Step-1**



a:5   b:9   C:12   d:13   e:16   f:45

**Step-2**

14 — a:5, b:9

**Step-3**

25 — C:12, d:13 ; 14 — a:5, b:9 ; e:16 ; f:45

**Step-4**

30 — [14 — a:5, b:9 ; e:16] ; 25 — C:12, d:13 ; f:45

**Step-5**

(tree: 55 → 30, f:45; 30 → 14, e:16; 14 → a:5, b:9; 25 → C:12, d:13)

Box:
$$\text{Compressed bits by Huffman encoding} = n \sum P(x)d(x)$$
$$B(T) = 100(2\cdot24)$$
$$= 224 \text{ bits}$$

**Step-6**

(tree: 100 → 0: f:45, 1: 55; 55 → 0: 25, 1: 30; 25 → 0: C:12, 1: d:13; 30 → 0: 14, 1: e:16; 14 → 0: a:5, 1: b:9)

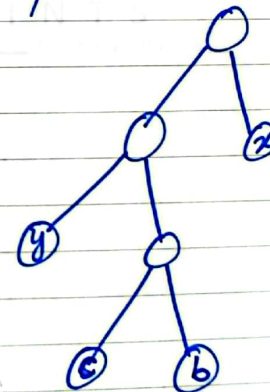| Codes | bits | |
|---|---|---|
| a ⇒ 1100 | 4 × 5 = 20 | |
| b ⇒ 1101 | 4 × 9 = 36 | Total bits |
| c ⇒ 100 | 3 × 12 = 36 | |
| d ⇒ 101 | 3 × 13 = 39 | = 224 |
| e ⇒ 111 | 3 × 16 = 48 | |
| f ⇒ 0 | 1 × 45 = 45 | |

---

# Proof of Correctness

## Greedy choic Property :-

Let two characters 'x' and 'y' with smallest probabilities, Then there is optimal code tree in which 'x' and 'y' are siblings at maximum depth.

## Proof

Let 'T' be any optimal prefix code tree with two siblings 'b' and 'c' at max. depth.



As

$$P(b) \geq P(x)$$
$$P(b) - P(x) \geq 0 \quad \text{———①}$$

&

$$d(b) \geq d(x)$$
$$d(b) - d(x) \geq 0 \quad \text{———②}$$

Multiplying ① and ②

$$\left[P(b) - P(x)\right]\left[d(b) - d(x)\right] \geq 0 \quad \text{——③}$$

**Now** try to change our own Tree "T" to prefix code tree of Huffman for that, swap 'x' with 'b' to be 'x' (with Lowest probability/freq.) at more depth.

Now tree becomes "T'"

$$B(T') = B(T) - P(x)d(x) + P(b)d(x)$$
$$\qquad\qquad - P(b)d(b) + P(x)d(b)$$

$$= B(T) - \left[P(b) - P(x)\right]\left[d(b) - d(x)\right]$$

$$B(T') = B(T) - (\text{+ve factor})$$

Implying that

T' is optimal.

By swapping y and c in T, we'll get

T". So using previous argument we can say T" is also optimal.

So This shows that greedy choice made by Huffman also is proper one

## Optimal Substructure :-
(Proof by Induction)

## Base Case :-

For base case,

$$n = 1 \qquad \text{or} \qquad n = 2$$

| Tree consist of single leaf node (just 1 character) | Tree consist of two nodes |
|---|---|
| So that 1 character can be coded as '0' or '1' | code assigned to 1st character = 0 2nd character = 1 |
| So code length for that single character is 1 bit (which is optimal) | Again code length for them is 1 bit. (Obviously Optimal for compression) |

# Inductive hypothesis :-

"Let Huffman tree for "n-1"
characters is optimal."

# Inductive Case :-

We want to prove that
Huffman tree is also optimal for
exactly "n" characters.

if there are "n" characters in "T"

Now

Remove two characters (say 'x' and 'y')
with lowest propabilities at max depth;
and replace with character z
Such that

$$P(z) = P(x) + P(y)$$

Now

there are "n-1" characters
in tree (say "T'")

we can convert this tree T' to
the tree of 'n' characters T by
replacing z with nodes x and y

---

Now cost of tree T is

$$B(T) = B(T') - P(z)d(z) +$$
$$P(x)[d(z)+1] + P(y)[d(z)+1]$$

$$= B(T') - [P(x)+P(y)]d(z)$$
$$+ (d(z)+1)[P(x)+P(y)]$$

$$= B(T') + [P(x)+P(y)][-d(z)+d(z)]$$
$$+1]$$

$$B(T) = B(T') + \underbrace{[P(x)+P(y)]}_{\text{(no depth involved)}}$$

cost changes but change depends
in no way on ~~struct~~ depth of
tree

As our inductive hypothesis say that
tree (T) built optimal for n-1 characters [B(T')]

So
tree for n characters will also be
optimal (according to above equation)