## Huffman Coding

→ Those characters which have high freq. should have less storage space → Motive of Huffmann coding

Eg "a" in "aaabcd", so here 'a' shouldn't take 8 Bits / 1 Byte

This will never happen →

a → 00    g → 001

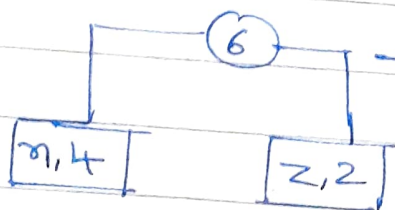Code = 00100  (but for ag you wanted 0001)
         ‿  ‿
         a   gx

Not Possible since these are Prefix free cod
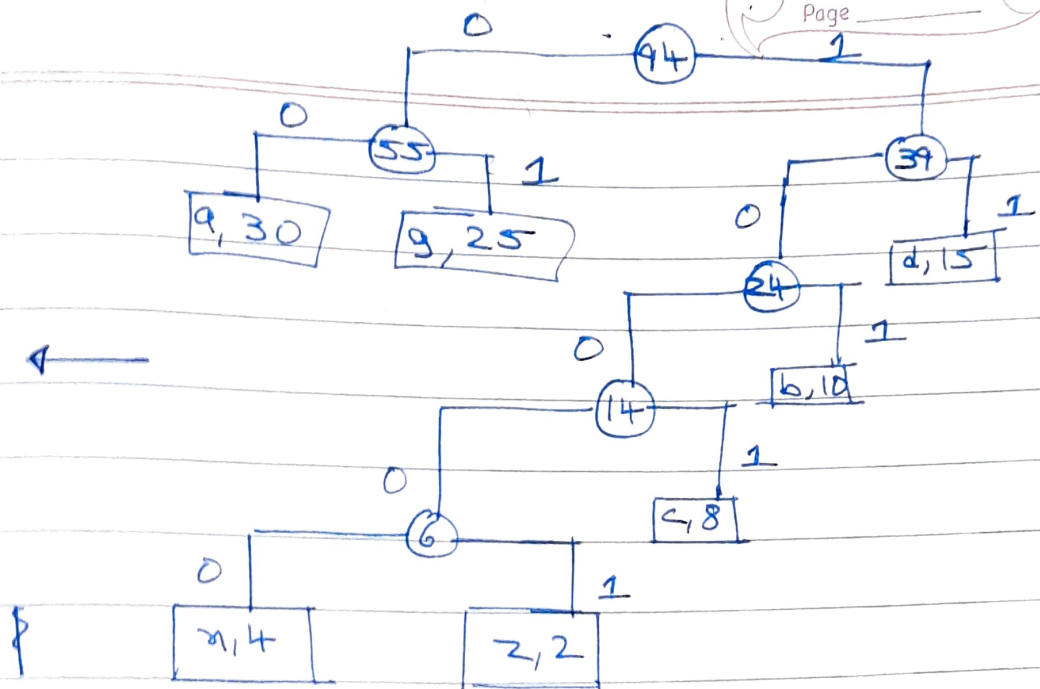
<u>Eg</u>

|  | Frequency |
|---|---|
| a | 30 |
| b | 10 |
| c | 8 |
| d | 15 |
| g | 25 |
| x | 4 |
| z | 2 |

Pick those characters which have Min. frequency



⑥ → Create a Node with freq. = Sum of freq. of these two and remove x,4 and z,2 from the frequency dictionary

x,4      z,2

94

0      1

55

0        1

| a, 30 | | g, 25 |

39

0      1

24

| d, 15 |

0      1

14

| b, 10 |

0      1

6

| c, 8 |

0      1

| x, 4 | | z, 2 |

This is a Tree ←

Now, to find code for each character,
→ For every Left Node add 0
→ For every Right Node add 1     ①→

∴ Code ⟹   a → 00

         g → 01

         d → 11

         b → 101

         c → 1001

         x → 10000

         z → 10001

5 Bytes
(1 Byte for each char) ①m

⟹ abadg →

code =
    0010100110|
    ∧

1 Byte    1 Byte ] → 2
(8 Bits)         Bytes

[Space reduced compared to ①]

So we need to stop adding 0/1 to code once we reach a Leaf node or Root Becomes None

→ To store freq of each character, we use Hash Map.

→ At each step, you want those 2 characters which have Min. frequency ⟶ We use Min Heap

→ Tree like Struct. ⟶ We use Tree

→ To store Code of each character ⟶ We use Hash Map (of Character, Bits)

- After Encoding ⟹

  a~~bc~~a
    ↓

  1011101~~1100~~
    ↓
  8 Bits

  abca
    ↓

  1011 1011 100
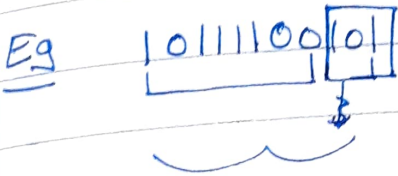  ⌊___⌋⌊___⌋⌊_⌋

  8 Bits   It'll convert this also into 8 Bits by     adding zeroes

  i.e. 10000000 ⟶ Aisa ho jayga

  ∴ While decoding, we don't know how many zeroes were present in the original code

length of

Make sure, BITS are Multiple of 8     text
To achieve this, you need to Pad this encoded.

How much part we have padded we will come
to know through the first 8 bits

Eg     | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 |

Total length = 11

$11 \% 8 = 3 =$ length of ☐ portion

$8 - (11 \% 8) = 5 =$ No. of zeroes to be added
after ☐

- Padding means adding zeroes at the end

- 1st 8 Bits Tell us how much we have padded

Eg If you have padded 3 zeroes, you'll store
00000011 in 1st 8 Bits
     Binary form of 3

If ___ 4 zeroes, ___ 00000100 ___
                              Binary form of 4

☆ ← (1st Pg)

## Problem Solving Tips

① Try calling function first and then doing what you wanted to do

☆ Data Structures req. for Huffman Coding:
① To store freq. of each char → Hash Map of chars and integers
② To Build a Tree → Tree
③ I want chars with Min. freq. → Min Heap
④ To store Code of each char → Hash Map of chars and Bits

Eg  a → 00

b → 101

• Decompression means Decoding
Eg  00 means 'a', 101 means 'b'
Problem →

a → 00    g → 001

String = 00100   OR  00100
        ‿ ‿          ‿ ‿
        a  x         g  x   This will never

happen (Codes are Prefix free) i.e. if a is 00 no code of any other char (say g will start from 00)

d)

Padding → 1 0 1 1 1 1 1 0 0 1 0 1

⟨8 Bits⟩  Add 5 '0's' to make 8 Bits

⟹ 1 0 1 1 1 1 1 0 0 1 0 1 0 0 0 0 0

⟨8 Bits⟩  ⟨8 Bits⟩

~~For me~~ I should know that I have added 5 '0s' to pad the text, so this padded info, I'll store in 1$^{st}$ 8 Bits. Here, I have padded 5 zeroes, Binary form of 5=101 I want to store it as a Binary No. consisting of 8 Bits

∴ 0 0 0 0 0 1 0 1  1 0 1 1 1 1 0 0  1 0 1 0 0 0 0 0

Padded_info  8 Bits  8 Bits

( Binary form of 5)

To Remove Padding,

if original text = [10011100][10001]

⟹ Padded text = [10011100] [10001 000]

Padding (e=3)

∴ org.text = Padded text [ :-1 * e]