

(12) File Zipper (Huffman Encoding) 25/02/21

Just like WinRAR we'll create a project to which could compress files (for now only text(.txt) files) and it'll be a lossless compression that is there'll be no data loss.

In our project we'll display 2 things, first will be the encoded data or the compressed data and second will be a Compression Ratio (If its greater than 1 then compression is lossless or successful! and higher it goes the better compression gets)

◆ Huffman Coding

> It's a lossless data compression algorithm

> We assign variable-length codes to input characters, length of which depends on frequency of chars:

ex. $\frac{aaaaa}{5} \frac{bb}{45} \rightarrow$

	Freq.	Code
a	5	1100
b	45	0

> The variable-length codes assigned to input characters are prefix codes.

$\{0, 11\}$
Prefix Codes

①

$\{0, 1, 11\}$
Non prefix Codes

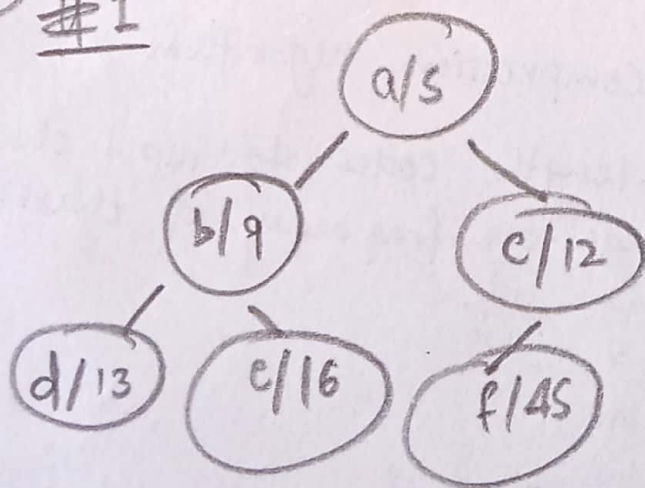
this makes sure our data is uniquely decodable

• Algo to create Huffman Tree

- ① Create a leaf node for each unique character and build a min heap of all leaf nodes.
- ② Extract two nodes with the min. freq. from the min. heap.
- ③ Create a new internal node with freq. equal to the sum of the 2 nodes frequencies. Make the first node as its left child and the other extracted node as its right child. Add this node to the min heap.
- ④ Repeat steps #2 and #3 until the heap contains only one node. The remaining node is the root node and the tree is complete.

ex

#1



Min Heap (MH)

char	freq
a	5
b	9
c	12
d	13
e	16
f	45

#2

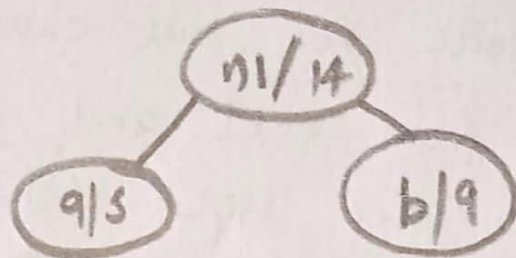
a/s

&

b/9

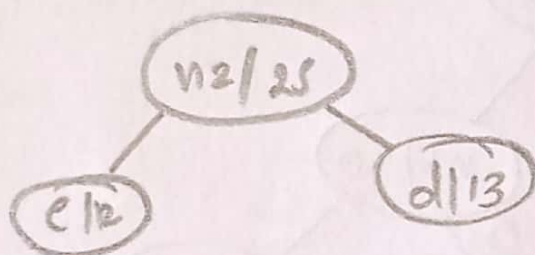
②

3



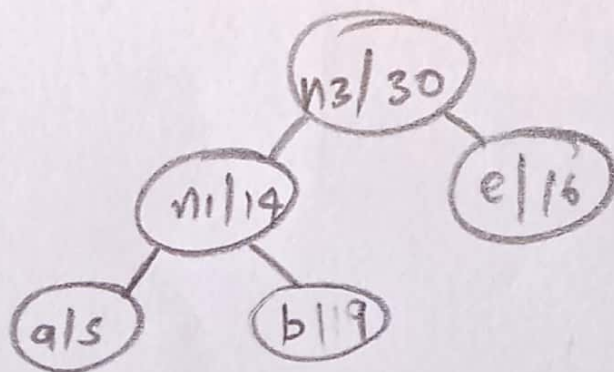
char	freq	MM
n1	14	
c	12	
d	13	
e	16	
f	45	

#4 → #2 & #3



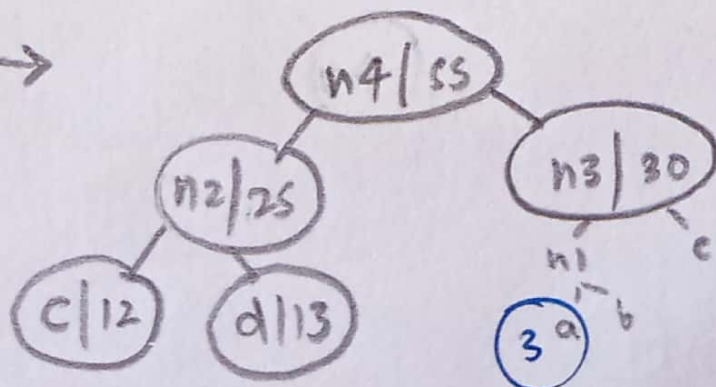
char	freq	MM
n1	14	
n2	25	
e	16	
f	45	

→ #2 & #3

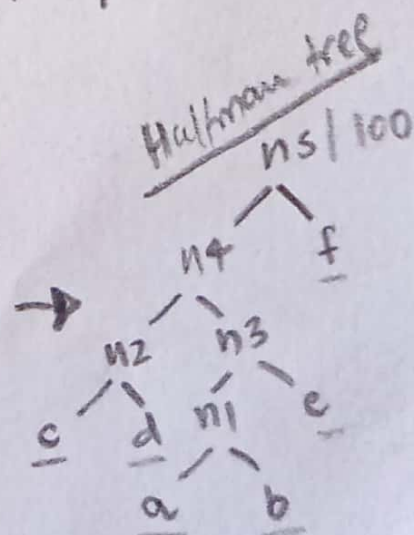


char	freq	MM
n2	25	
n3	30	
f	45	

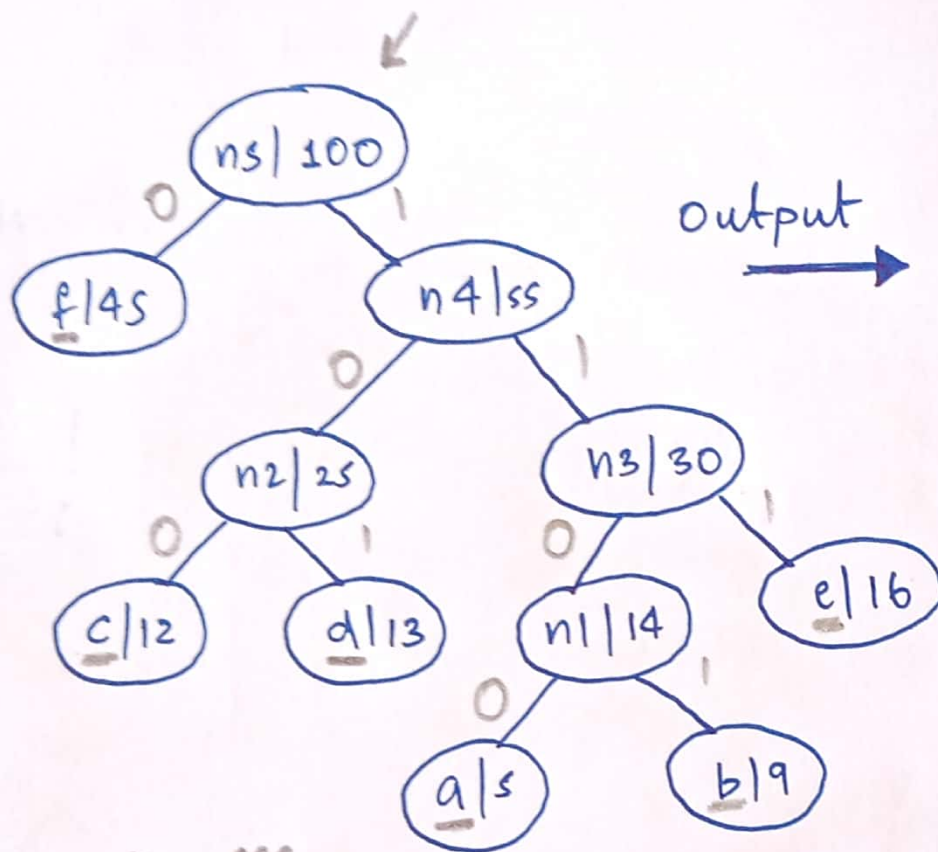
→



c	f
n4	55
f	45



Now we'll generate code for our chars through this tree by traversing from root and assign '0' to left child and '1' to the right child, and once we encounter a leaf node we print the code.



char	Code
a	1100
b	1101
c	100
d	101
e	111
f	0

As we can see that higher the freq. of a char lower is its code length.

$$\therefore \text{len of Code} \propto \frac{1}{\text{freq. of char.}}$$

take one more ex.

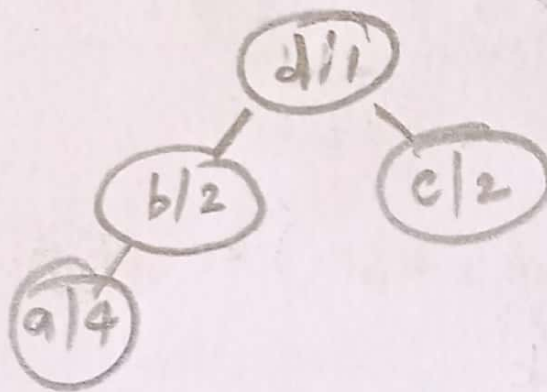
Input → a b a c d a b a c

• first we'll create a freq. arr. for all unique chars.

we get → a-4
b-2
c-2
d-1

• Now we'll create leaf node (of char & freq) for all unique chars and build min heap with 'em.

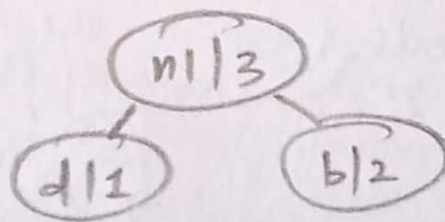
#1



char	freq
d	1
b	2
c	2
a	4

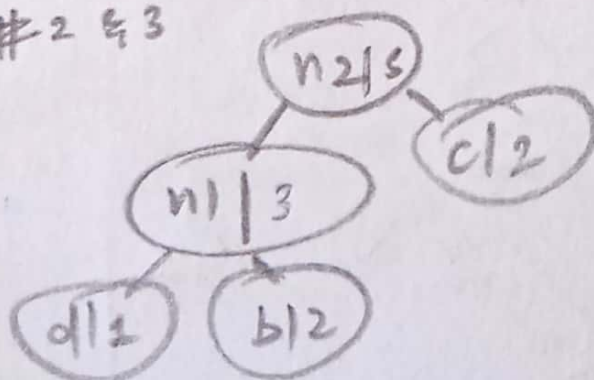
• Now take 2 nodes with min freq. & make a node of them.

#2 & 3



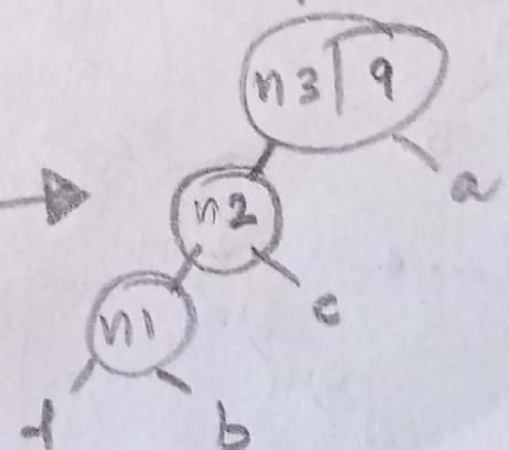
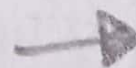
char	freq
n1	3
c	2
a	4

• #2 & 3

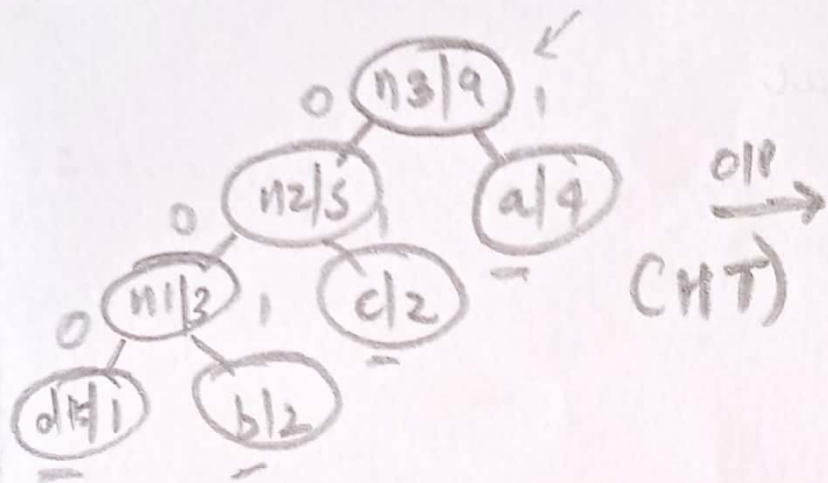


c	f
n2	5
a	4

(5)



• Acc. we'll write code for it:



char	code
d	0 0 0
b	0 0 1
c	0 1
a	1

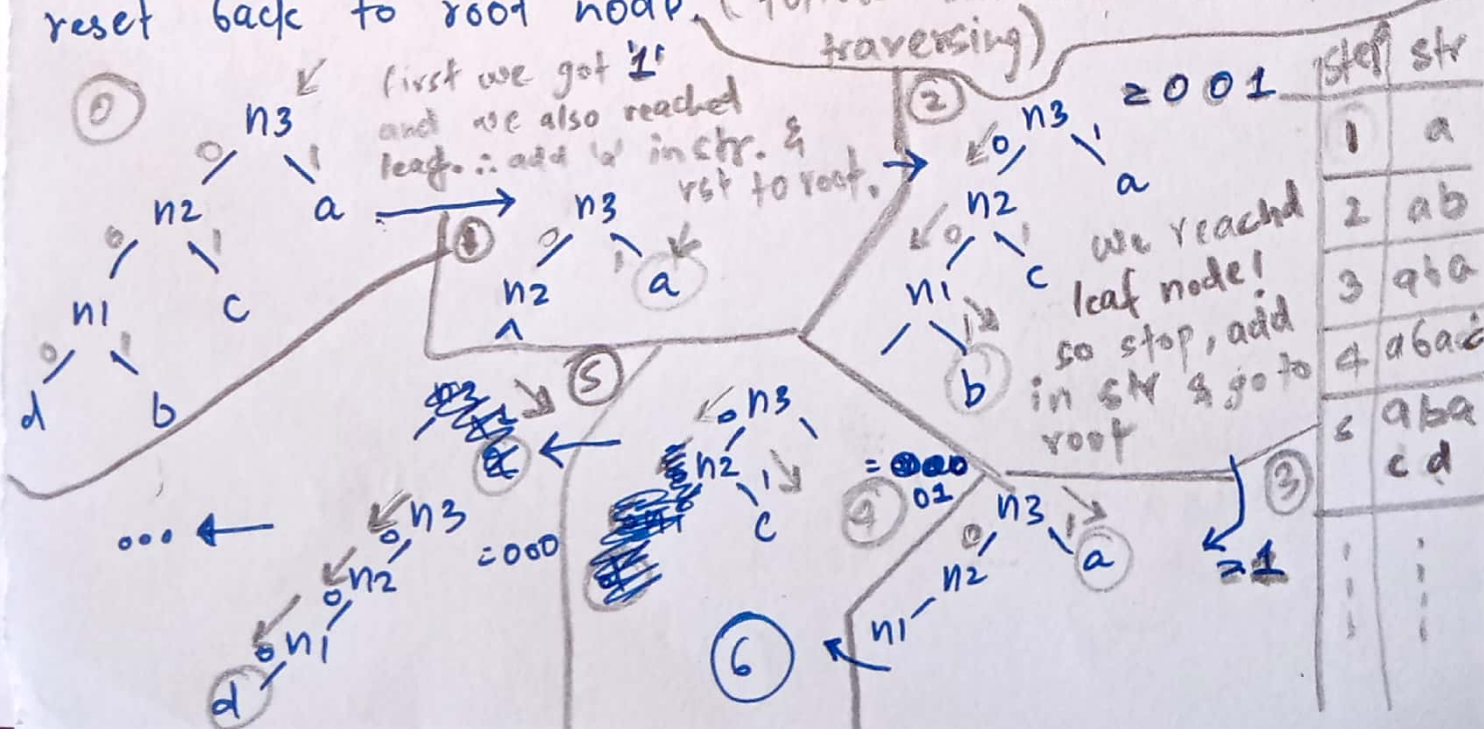
∴ with the help of our O/P we can generate the code:

abac dabac $\xrightarrow{\text{encoded}}$ 1 001 1 01 000 1 001 1 01
17 bits

∴ After encoding we reduced 72 bits (9 chars x 8 bits) to 17 bits!

• Now let decode it back:

- traverse from root. create an empty string. if we reach a leaf node then add its char to string & reset back to root node. (follow encoded data for traversing)



```
import { BinaryHeap } from './heap.js';
export { HuffmanCoder }
```

```
class HuffmanCoder {
```

```
  stringify(node) {
```

```
    if (typeof(node[1]) == "string") {
      return '\ ' + node[1];
    }
```

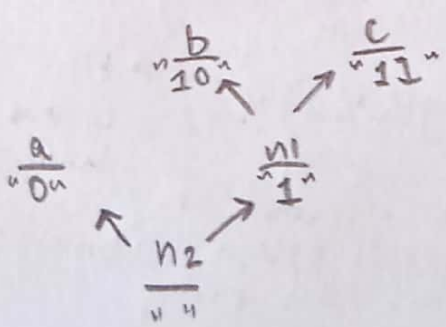
```
    return '0' + this.stringify(node[1][0]) +
      '1' + this.stringify(node[1][1]);
  }
```

```
  getMappings(node, path) { → follows DFS
```

```
    if (typeof(node[1]) == "string") {
      this.mappings[node[1]] = path;
      return;
    }
```

```
    this.getMappings(node[1][0], path + "0");
    this.getMappings(node[1][1], path + "1");
  }
```

To leaf nodes
honge unke 2nd
para charctr hoga
and in js
chars and str
are reprsntd
similarly



7

5

encode(data) {

→ Jo string ko encode krna hota hai ex: 'abc'

this is the ← { this.heap = new BinaryHeap();

max heap

const mp = new Map();

for (let i=0; i<data.length; i++) {

if (data[i] in mp) {

mp[data[i]] = mp[data[i]] + 1;

}

else {

mp[data[i]] = 1;

}

}

for (const key in mp) {

this.heap.insert([-mp[key], key]);

}

while (this.heap.size() > 1) {

const node1 = this.heap.extractMax();

const node2 =

"

;

const node = [node1[0] + node2[0], [node1, node2]];

}

③

④ { const huffman_encoder = this.heap.extractMax();

⑤ { this.mappings = {};

this.getMappings(huffman_encoder, " ");

let binary_string = "";

⑥

for (let i=0; i<data.length; i++) {

binary_string = binary_string + this.mappings[data[i]];

}

we also need to store it as while decoding
apko original string chahiye hogi to decode

let rem = (8 - binary_string.length % 8) % 8;

let padding = "";

⑦

for (let i=0; i<rem; i++) {

padding += "0";

binary_string += padding;

⑧

ye utne 0's
string format me
store krlega thr binary-string
me add krdeg

kitne extra 0 bits
chahiye & bit
banane ke liye

let result = "";

for (let i=0; i < binary-string.length; i+=8) {

let num = 0;

for (let j=0; j < 8; j++) {

num = num * 2 + (binary-string[i+j] - "0");

}

result = result + String.fromCharCode(num);

}

9 { let final_res = this.stringify(huffman_encoder) + '\n' +
sem + '\n' + result;

10 { let info = "Compression ratio: " + data.length / final_res.length;
info = "Compression ratio complete & file sent for download" + info;

11 { return [final_res, this.display(huffman_encoder, false), info];
}

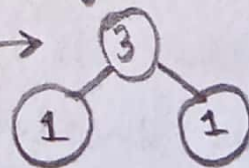
Code explanation for encode fn() :

1 We made an frequency hashmap for each character, we can also create an array of size 256.

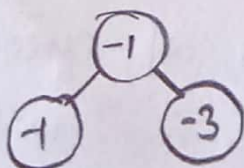
2 Now we add all the chars along with their freq. in heap. As it's a max. heap we add -freq. which gives us a min heap! ^{ex. (3, 'a')}

ex We had map = { {a:3},
{b:1},
{c:1} }

so if apn '-' nhi lagate to heap
esa banta →



It's just a
hack to get
min heap.

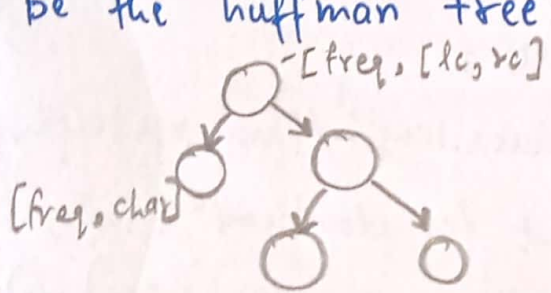


but as we use
min heap and pick top 2 nodes
with min freq. we add '-' in
front of every freq & we get:

9

③ Here while we have more than 1 nodes in heap we take top 2 nodes ^(n_1 & n_2) with min freqs and create a new internal node with n_1 & n_2 's freq sum as 1st para and an array of left & right child as 2nd para. Then we finally insert the newly created node back in the heap.

④ Now finally we'll extract the root node which'll be the huffman tree. ^(HT)



Structure of HT:

Internal node

all the non-leaf nodes will have a node structure of freq. and an arr of size 2 with left child & right child. ($[freq, [-, -]]$) and all leaf nodes will have freq & char only. ($[freq, char]$)

⑤ Through this we get the code for all the leaf nodes, in mapping. ^{this gave us} mappings =

$\{ \begin{array}{l} a: 0 \\ b: 10 \\ c: 11 \end{array} \}$

⑥ This creates the encoded code by apndⁿ all the codes of chars $\rightarrow 0101100$

⑦ & ⑧ Initially we had data = "abcaa" which was of $8 \times 5 = 40$ bits, Now although we have encoded it, it still is in the form of char which is of $7 \times 5 = 35$ bits! to reduce it: ~~as~~ ~~about~~ we'll convert them to ASCII chars! Now as ~~of~~ one char is of 8 bits we'll take ~~start~~ 8 chars from our encoded data and convert to ASCII charctr then we'll do with the remaining chars of encoded data, if the remaining chars are less than 8 we'll add remaining zeroes and then convert to ASCII this'll give us 0101100 & 0 = 8 we'll convrt this to it's ASCII & save it.

will reduce our bits from 40 to 8! & from 5 chars 1 char.

⑧ specifically convert our 8 bit codes to ASCII chars and ~~add~~ append to result. finally our result will have all 8 bit ASCII chars.

* After encryption our data abcaa will become X as 01011000's ASCII is X. and X doesn't mean anything in

⑨ stringf func gets string var of our Huffman tree. n apr 'v' islye add kote str-func me kyoki agr apne data me 0 or 1 kote to differentiate unhi kote pate path (ie 0s & 1s) se islye we use 'x'.

⑩ To decode the code apr final-res banate hai which apart from result extra info store karta hai like remaining bits and string var. of our tree to finally apne bits 40 se lagbhg 16 to gaye which is good!

* Note → In some cases our compression might not work. when our data is lrs like abc (24 bits) our final-res might come (50 bits). so to get meaningful or get data which has more than 10k or grtr chars.

⑩
$$\text{Compression ratio} = \frac{\text{original data len}}{\text{new}} \Rightarrow \frac{5}{1} = 5$$

this tells how good compressⁿ was.

⑪ Finally we return all the informatⁿ like final-res, info and also call display fn to display our tree.

- Remaining code of Huffman Codes Class

destringify (data) {

let node = [];

if (data[this.ind] == '\') {

this.ind++;

node.push(data[this.ind]);

this.ind++;

return node;

}

this.ind++;

let left = this.destringify(data);

node.push(left);

this.ind++;

let right = this.destringify(data);

node.push(right);

return node;

}

decode (data) {

data = data.split('\n'); } ①

if (data.length == 4) {

data[0] = data[0] + '\n' + data[1];

data[1] = data[2];

data[2] = data[3];

data.pop();

}

③ { this.ind = 0;

const huffman_decoder = this.destringify(data[0]);
const text = data[2];

let binary_string = "";

for (let i=0; i<text.length; i++) {

let num = text[i].charCodeAt(0);

let bin = "";

for (let j=0; j<8; j++) {

bin = num % 2 + bin;

num = Math.floor(num/2);

}

binary_string += bin;

}

to make our
binary string of 8
bits we added some
padding. So to remove
them we use data[1]
and get our original
binary string back.

④ { binary_string = binary_string.substring(0, bin_str.length - data[1])

console.log(binary_string.length);

let res = "";

let node = huffman_decoder;

⑤ { for (let i=0; i<binary_string.length; i++) {

if (binary_string[i] === '0') {

node = node[0];

else {

node = node[1];

}

if (typeof (node[0]) === "string") {

res += node[0];

node = huffman_decoder;

}

}

⑥ { let info = "decompression";
return [res, this.display(huffman_decoder, true), info];

• Explanation

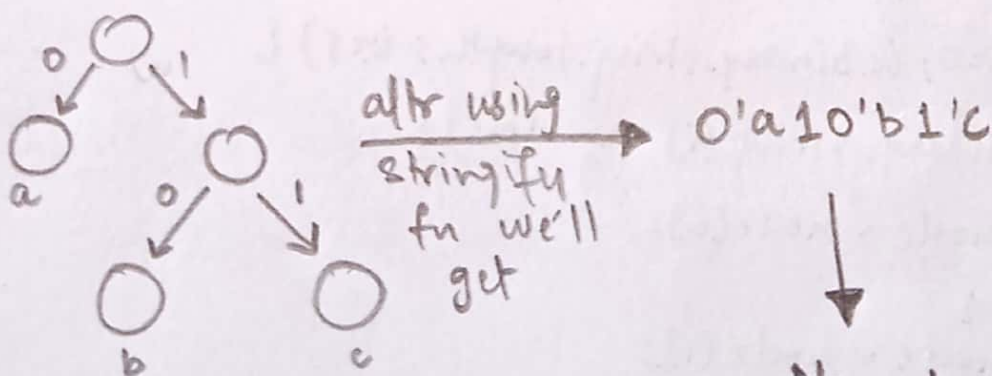
From the encode func we receive data in the form of "^{HT}Tree \n ^{padding}Rem \n Encoded data".

① ∴ To retrieve these information we use split fn which splits data on the basis of "\n". That's how we get 3 diff. info in data.

② Handles Special Case. If we have '\n' (new line) in our HT then we do this.

eg) In our tree we get '\n' → Tree
1st para '\n'
2nd para, etc
∴ Instead of 2 '\n' we get 3 '\n' and string breaks into 4 parts.

③ Now we have string vs of HT in our data[0], we need to get our original HT to decode the data. Let's take an ex to understand:-



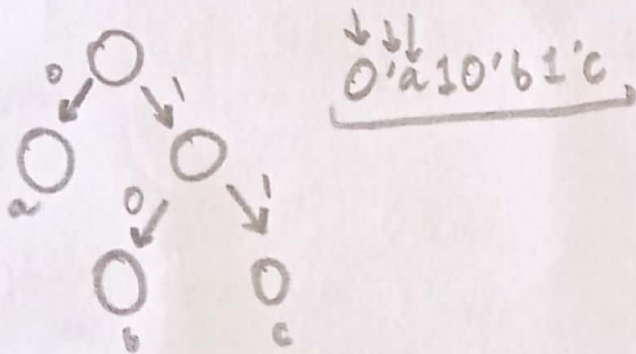
Now to deserialize it or to get the original HT we'll handle 2 case:

③i For leaf node

③ii For internal node

(3ii) Every internal node will have left and right child. (3ii) will go on these nodes. By marking "ind" we ensure that we catch only leaf nodes, backi for left and right children (ie for 0 & 1) we just keep traversing.

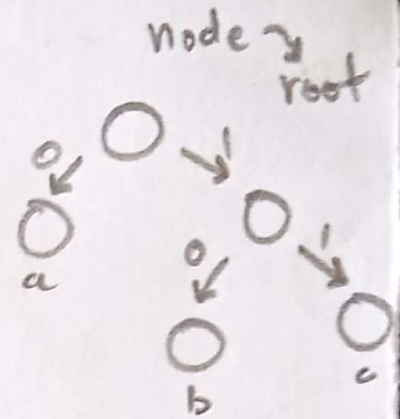
(3i) As soon as we get '0' we can infer that we've got to leaf node. After reaching we increment "ind" to reach charctr then push it in node then again start to read next char.



(4) Now we've to get our binary string from encoded text (ex) "!, " jaise code se → 1001101... so that we could decode it back! suppose our encoded text was → "!, " now this was created via binary string, now we've to again get back our binary string via encoded text. let's assume for "!" ASCII code is 1000110, first we get "!" number (that ranges from 0-256) then we get its binary string and append it in "binary-string", similarly we do with " ,".

④ Now here we'll use our binary-string after removing the padding, Huffman Tree (HT) to get back the original data

from "!,," $\xrightarrow[\text{after rem'n padding}]{\text{we got}}$ 0110010
 \uparrow binary-string



we'll traverse the binary string, first we start from root node. then we got "0" so we update our node to left child node, like this we keep on going (for 0 we go left & for 1 right).

And when we get leaf node, we add its char to result ("res") and reset our node to root node.

⑥ finally we get our decoded data in "res".
In end return.