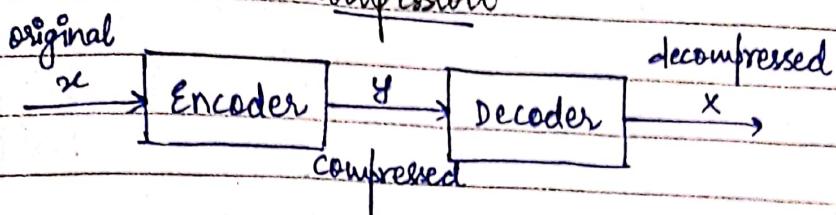


24/01/19

UNIT : 1Compression➤ Classification

Compression is divided in 2 parts :-

(i) Lossless compression

- Decoded data = original data
 - Compression ratio < lossy compression
- Eg. LZ, JBIG

(ii) lossy compression

- Decoded data ~ original data
- Compression ratio > lossless compression
- Used where small errors are allowed.

Eg. MPEG, JPEG.

➤ Entropy

It's the measurement of info content for a set of message s with $P(E)$ $P(S)$.

Self info of s is :-

$$i(s) = \log \frac{1}{P(s)}$$

• conditional entropy

$$H(S) = \sum_{s \in S} P(s) \log \frac{1}{P(s)}$$

Q. $P(S) = \{0.25, 0.25, 0.25, 0.125, 0.125\}$

Find $H(S)$

$$H(S) = 3 \left[0.25 \log_2 \frac{1}{0.25} \right] + 2 \left[0.125 \log_2 \frac{1}{0.125} \right]$$
$$= 3(0.25 \log_2 4) + 2(0.125 \log_2 8)$$
$$= 2.25$$

Q. $P(S) = \{0.5, 0.125, 0.125, 0.125, 0.125\}$

Find $H(S)$

$$H(S) = 1 \left[0.5 \log_2 \frac{1}{0.5} \right] + 4 \left[0.125 \log_2 \frac{1}{0.125} \right]$$
$$= 2$$

Q. $P(S) = \{0.75, 0.0625, 0.0625, 0.0625, 0.0625\}$

Find $H(S)$

$$H(S) = 1 \left[0.75 \times \log_2 \frac{1}{0.75} \right] + 4 \left[0.0625 \log_2 \frac{1}{0.0625} \right]$$
$$= \sqrt{2}$$
$$= 1.414$$

► Measures of Performance

Compression algo can be evaluated in a no. of diff' ways. We can measure the relative complexity of the algorithm in terms of compression ratio.

Compression ratio = $\frac{\text{bits required to represent data before compression}}{\text{bits required to represent data after compression}}$

Q. Suppose storing an image of size 256×256 pixel requires 65538 bytes. After compression it requires 16384 bytes.
 Compression ratio = $\frac{65538}{16384} \rightarrow 1:1$

Shanon Fano coding

- This is a source coding method to construct binary codes
- The codes should satisfy following conditions :-
 (i) No sequence of codes can be obtained from each other by adding more binary digit to shorter codes (prefix P(E))
 (ii) Tx of encoded message is efficient. 1 & 0 appear independently.

- Procedure

- Message are written in decreasing order of their P(E)s.
- The message set is divided into 2 subsets of equal or nearly equal P(E)s
- Code 0 is assign to each message in one subset above partition line.
- Code 1 is assign to each message in one subset below partition line.
- The procedure continue until each subset contains only one message.

Q. Apply shanon Fano coding procedure for following message

$$[x] = [x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8]$$

$$[P] = \left[\frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{16}, \frac{1}{16}, \frac{1}{4}, \frac{1}{16}, \frac{1}{8} \right]$$

take M = 2.

<u>Message</u>	<u>probabilities</u>	s_1	s_2	s_3	s_4
x_1	0.25	0	0		
x_6	0.25	0	1		
x_2	0.125	1	0	0	
x_8	0.125	1	0	1	
x_3	0.0625	1	1	0	0
x_4	0.0625	1	1	0	1
x_5	0.0625	1	1	1	0
x_7	0.0625	1	1	1	1

$$\text{efficiency} = \eta = \frac{H(X)}{\bar{L} \cdot \log_2 M}$$

$H(X)$ - Entropy

$$\bar{L} = \sum p_k n_k$$

$$\Rightarrow \bar{L} = \frac{1}{4} \times 2 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{8} \times 3 + \frac{1}{16} \times 4 + \frac{1}{16} \times 4 \\ + \frac{1}{16} \times 4 + \frac{1}{16} \times 4 = 2.75$$

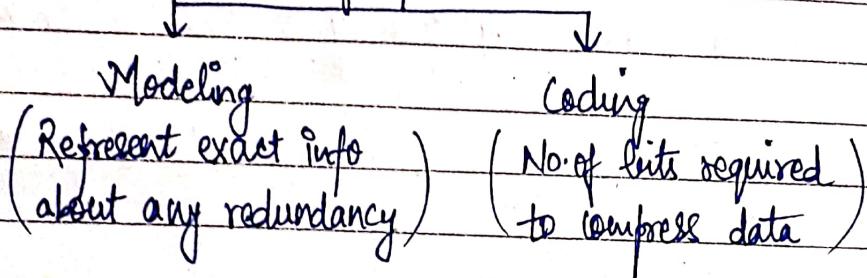
$$H(X) = 2 \left[0.25 \log \frac{1}{0.25} \right] + 2 \left[0.125 \log \frac{1}{0.125} \right] \\ + 4 \left[0.0625 \log \frac{1}{0.0625} \right] \\ = 1 + 0.75 + 1 = 2.75$$

$$\text{efficiency} = \eta = \frac{2.75}{2.75} = 1$$

► Modelling & coding

In data compression, when we reconstruct message, the decision whether a compression scheme is to be lossy or lossless, the exact compression scheme we use will depend on a no. of diffn factors. Some of the most important factors are characteristic of data that need to be compressed.

Data compression algorithm



• Models

(i) Physical model

If we know something about the physics of data generation, we can use that information to construct a model.

e.g. speech related application.

Knowledge about the physics of production can be used to construct a mathematical model for a sample speech model.

(ii) Probability model

The simplest statistical model for the source is to assume that each letter that is generated by the source is independent of every other letter & each occurrence of letter have same $P(E)$. This kind of model is known as $P(E)$ model. In this model, we consider only source.

(iii) Markov model

This is one of the most popular model. It represents dependency in data. The model is given by famous mathematician Markov. This is used in lossless compression. In this, we use discrete time frame $\{t_x\}$ & observation frame $\{X_n\}$.

The $P(E)$ of observation :-

$$P(X_n | X_{n-1}, \dots, X_{n-k}) = P(t_{X_n} | t_{X_{n-1}}, \dots, t_{X_{n-k}})$$

- In other words, knowledge of last ' k ' symbols is equivalent to the knowledge of entire past history of the process.
- The value taken by the set X_{n-1}, \dots, X_{n-k} are called the state of process.
- If the size of source alphabet is ' l ' then no. of states is ' l^k '.

$$X_n = P X_{n-1} + E_n$$

↓
Noise

- Q. Consider a binary image has a 2 type of pixel - black & white pixel. We know that the appearance of white pixel as the next observation depends on some extent whether the current pixel is white or black.

$$\begin{array}{ll} p(w/b) & S_w \rightarrow \text{current pixel is black} \\ p(b/w) & S_f \end{array}$$

$$\text{entropy } H = \sum_{i=1}^M p(S_i) H(S_i)$$

$$H(S_w) = -p(b/w) \log p(b/w) - p(w/b) \log p(w/b).$$

31/01/17

► Coding

It means assignment of binary sequence to element of an alphabet.

e.g. letter codeword

a_1	0
a_2	00
a_3	1
a_4	11

Codeword: Individual no. of code known as codeword

Alphabet: It's a collection of symbols called letters

Rate of code: The average no. of bits per symbol is called as the rate of code.

- Uniquely decodable code

The average length of the code isn't the only important point in designing a good code.

Q. Suppose we have 4 letters :-

$$a_1 \quad P(a_1) = 1/2$$

$$a_2 \quad P(a_2) = 1/4$$

$$a_3 \quad P(a_3) = 1/8$$

$$a_4 \quad P(a_4) = 1/8$$

Code 1	Code 2	Code 3	Code 4
0	0	0	0
0	1	10	01
1	00	110	011
10	11	111	0111



- Average length (l) = $\sum p(a_i) n(a_i)$

$$l_1 = 0.5 \times 1 + 0.25 \times 1 + 0.125 \times 1 + 0.125 \times 2 \\ = 1.125 \text{ bits/symbols}$$

$$l_2 = 0.5 \times 1 + 0.25 \times 1 + 0.125 \times 2 + 0.125 \times 2 \\ = 1.25 \text{ bits/symbols}$$

$$l_3 = 0.5 \times 1 + 0.25 \times 2 + 0.125 \times 3 + 0.125 \times 3 \\ = 1.75 \text{ bits/symbols}$$

$$l_4 = 0.5 \times 1 + 0.25 \times 2 + 0.125 \times 3 + 0.125 \times 4 \\ = 1.875 \text{ bits/symbols.}$$

- Test for unique decodability

Suppose we have two binary codewords a, b . ' a ' is ' k ' bit long & ' b ' is ' n ' bit long & $k < n$. If the k bit of ' b ' are identical to ' a ' then ' a ' is called prefix of ' b '. $(n-k)$ bits of b called dangling suffix.

e.g. $a = 010$

$b = 010\underset{\substack{\text{prefix} \\ \text{dangling}}}{}11$
suffix

Steps for finding uniquely decodable codes

- Construct a list of all code words.
- Examine all pairs of codewords to check any codeword is prefix of another codeword.
- Whenever you find such pair, add dangling suffix to the list in previous iteration.
- Now repeat the procedure with current list if we get a dangling suffix that is a codeword in the list. Then the codeword isn't uniquely decodable

Eg. $a_1 = 0$
 $a_2 = 01$

$a_3 = 10$

$\{0, 01, 10\}$

$\{0, 01, 10, 1\}$

Not uniquely
decodable

Eg. $a_1 = 0 ; a_2 = 01 ; a_3 = 10 ; a_4 = 111$

$\{0, 01, 110, 111\}$

$\{0, 01, 110, 111, 1\}$

$\{0, 01, 110, 111, 1, 10\}$

$\{0, 01, 110, 111, 1, 10, 11\}$

Not uniquely
decodable

~~Prefix code~~

$\{0, 10, 110, 111\}$

$\{0, 10, 110, 111\}$

$\{0, 01, 011, 0111\}$

$\{0, 01, 011, 0111, 1, 11, 111\}$

0 01 110 111

a_1, a_2, a_3, a_4

a_1, a_2, a_3, a_4

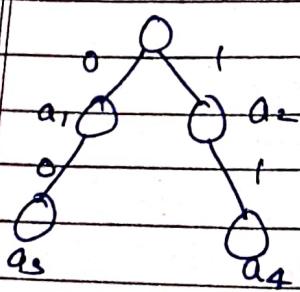
010219

Prefix code

A code in which no codeword is prefix to another codeword is called prefix code

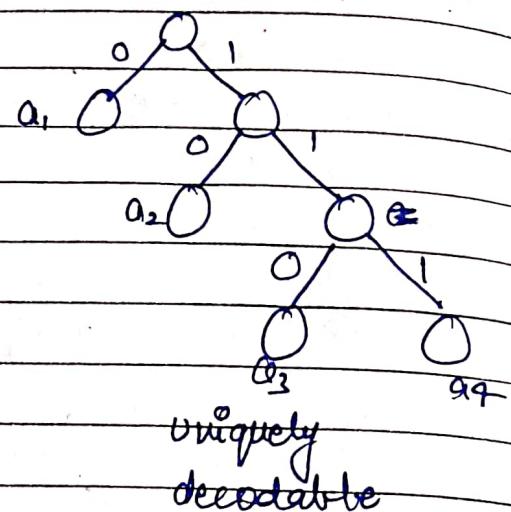
letter	code 1	code 2	code 3
a_1	0	0	0
a_2	1	10	01
a_3	00	110	011
a_4	11	111	0111

For code 1



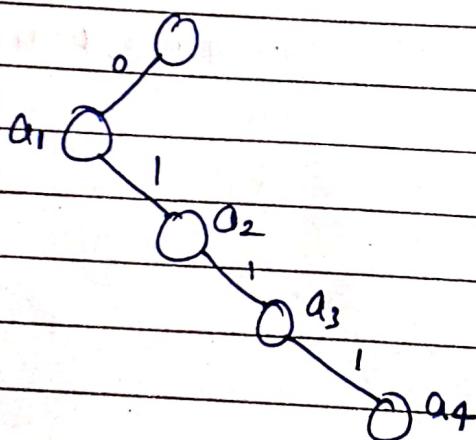
Not uniquely decodable

For code 2



Uniquely decodable

For code 3



Not uniquely decodable

Internal node: Node that gives rise to other node is called internal node.

External node: Node that doesn't give rise to other node.

→ Huffman Coding

This technique developed by David Huffman is a part of class assignment.

The code generated using this technique are called Huffman codes. These codes are prefix code & optimum for given model.

The Huffman procedure is based on two observation :-

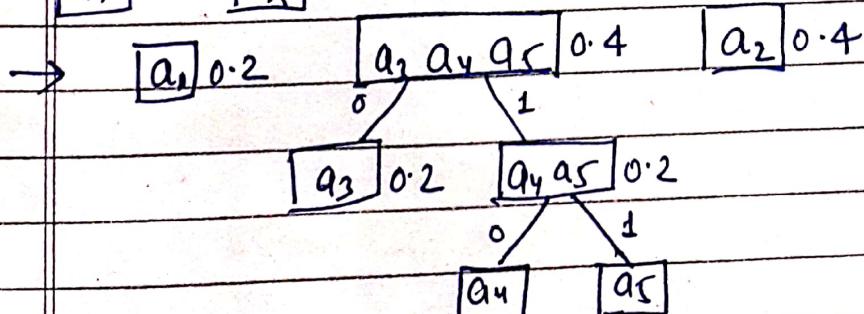
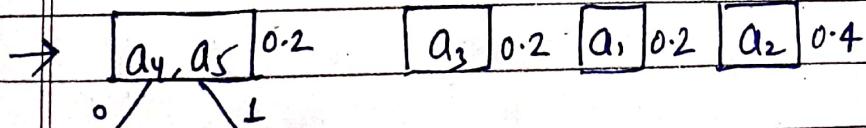
- In optimum code, symbols that occur more frequently will have shorter code words in comparison of that symbols occur less frequently.
- In an optimum code, the two symbols that occur least frequently will have the same length.

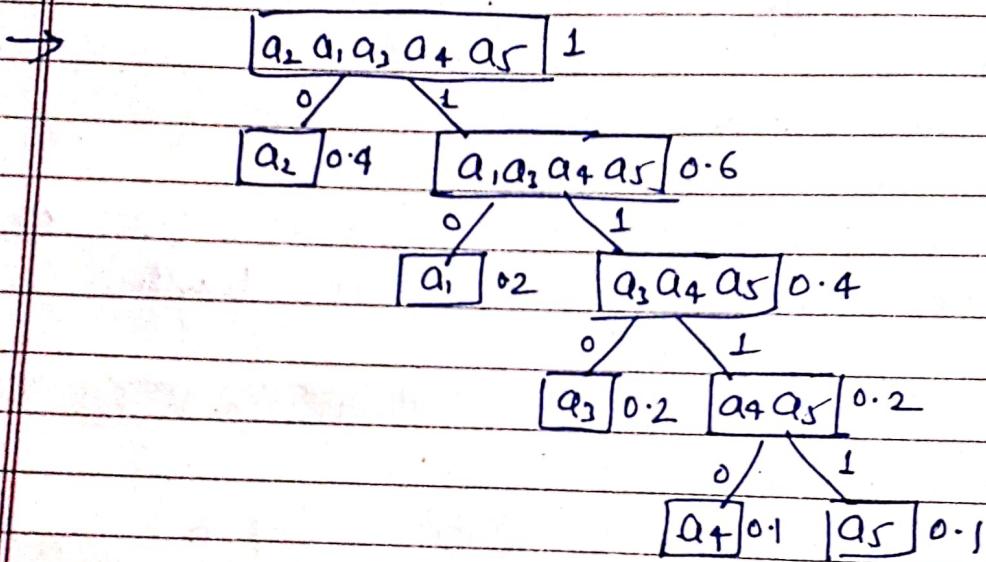
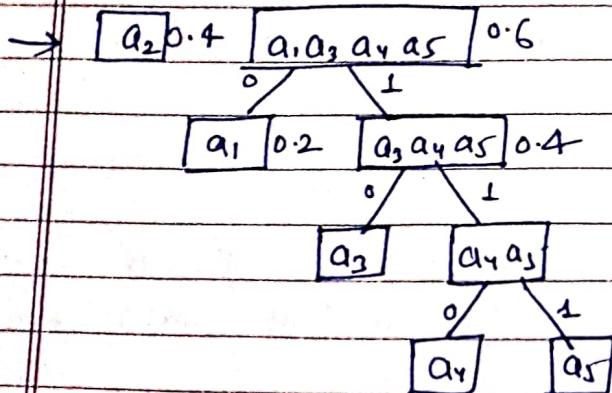
Q. Let $A = \{a_1, a_2, a_3, a_4, a_5\}$ with $P(a_1) = P(a_3) = 0.2$;

$P(a_4) = P(a_5) = 0.1$; $P(a_2) = 0.4$

Find Huffman code & calculate entropy.

a_1	0.2	a_2	0.4
a_2	0.4	\Rightarrow	a_1 0.2
a_3	0.2		a_3 0.2
a_4	0.1		a_4 0.1
a_5	0.1		a_5 0.1





$$\begin{aligned}
 \text{Entropy} &= 0.4 \times 1 + 0.2 \times 2 + 0.2 \times 3 + 0.1 \times 4 + 0.1 \times 4 \\
 &= 0.4 + 0.4 + 0.6 + 0.4 + 0.4 \\
 &= 2.2 \text{ bits / } \cancel{\text{signal symbol}}
 \end{aligned}$$

$$\text{Entropy} = \sum_{i=1}^n (\text{Code length} \times P(E_i))$$

8. A source emits letters from an alphabet
 $A = \{a_1, a_2, a_3, a_4\}$

$A = \{a_1, a_2, a_3, a_4, a_5\}$ with $p(E) \quad p(a_1) = 0.15, p(a_2) = 0.1$
 $p(a_3) = 0.26, p(a_4) = 0.15, p(a_5) = 0.5$

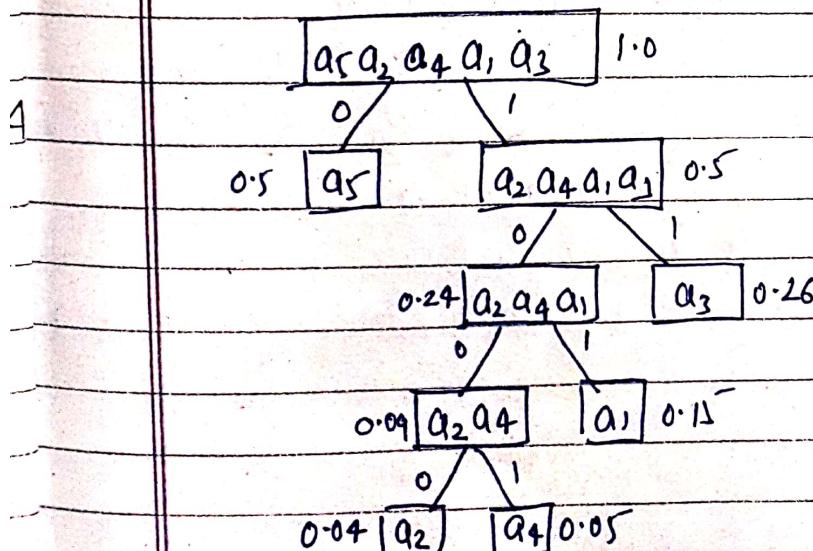
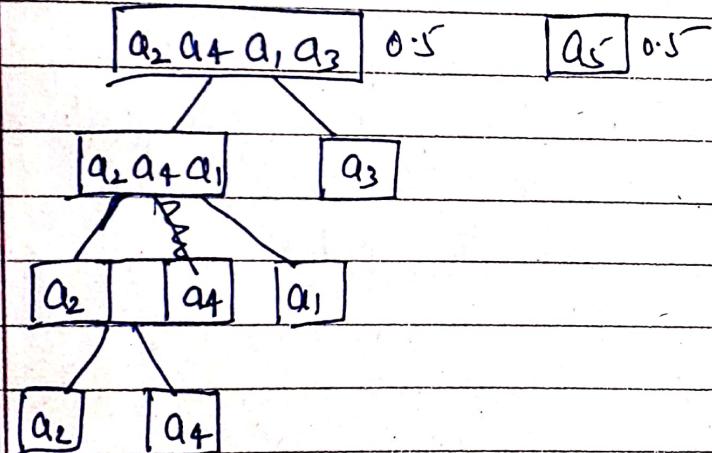
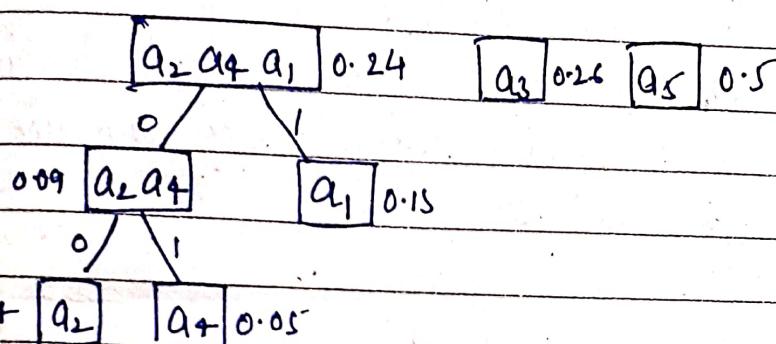
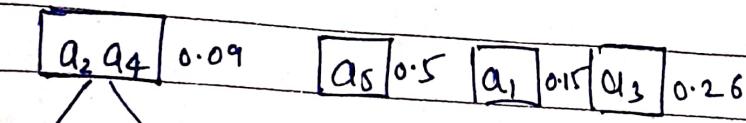
(i) calculate H.C.

(ii) find entropy

$a_1 \quad 0.15$
 $a_2 \quad 0.04$
 $a_3 \quad 0.26$
 $a_4 \quad 0.05$
 $a_5 \quad 0.5$

\Rightarrow

$a_5 \quad 0.5$
 $a_3 \quad 0.26$
 $a_1 \quad 0.15$
 $a_5 \quad 0.5$
 $a_4 \quad 0.05$
 $a_2 \quad 0.04$



$$\begin{aligned}
 \text{Entropy} &= 0.5 \times 1 + 0.26 \times 2 + 0.15 \times 3 + \cancel{0.05} \times 4 + 0.04 \times 4 \\
 &= 0.5 + 0.52 + 0.45 + 0.2 + 0.16 \\
 &= 1.02 + 0.45 + 0.36 \\
 &= 1.83
 \end{aligned}$$

05/02/19

Adaptive Huffman Coding

Huffman coding require knowledge of the $P(E)$'s of the source sequence. If this knowledge isn't available, we can use adaptive Huffman coding.

In the adaptive Huffman coding, neither transmitter nor receiver know anything about source sequence at the start of transmission.

In this procedure, transmitter & receiver consist of a single node which is known as NYT (Not Yet Transmitted) node & weight of this node is 0.

At the transmission progress, nodes corresponding to symbols transmitted will be added to tree & tree is reconfigured using update procedure.

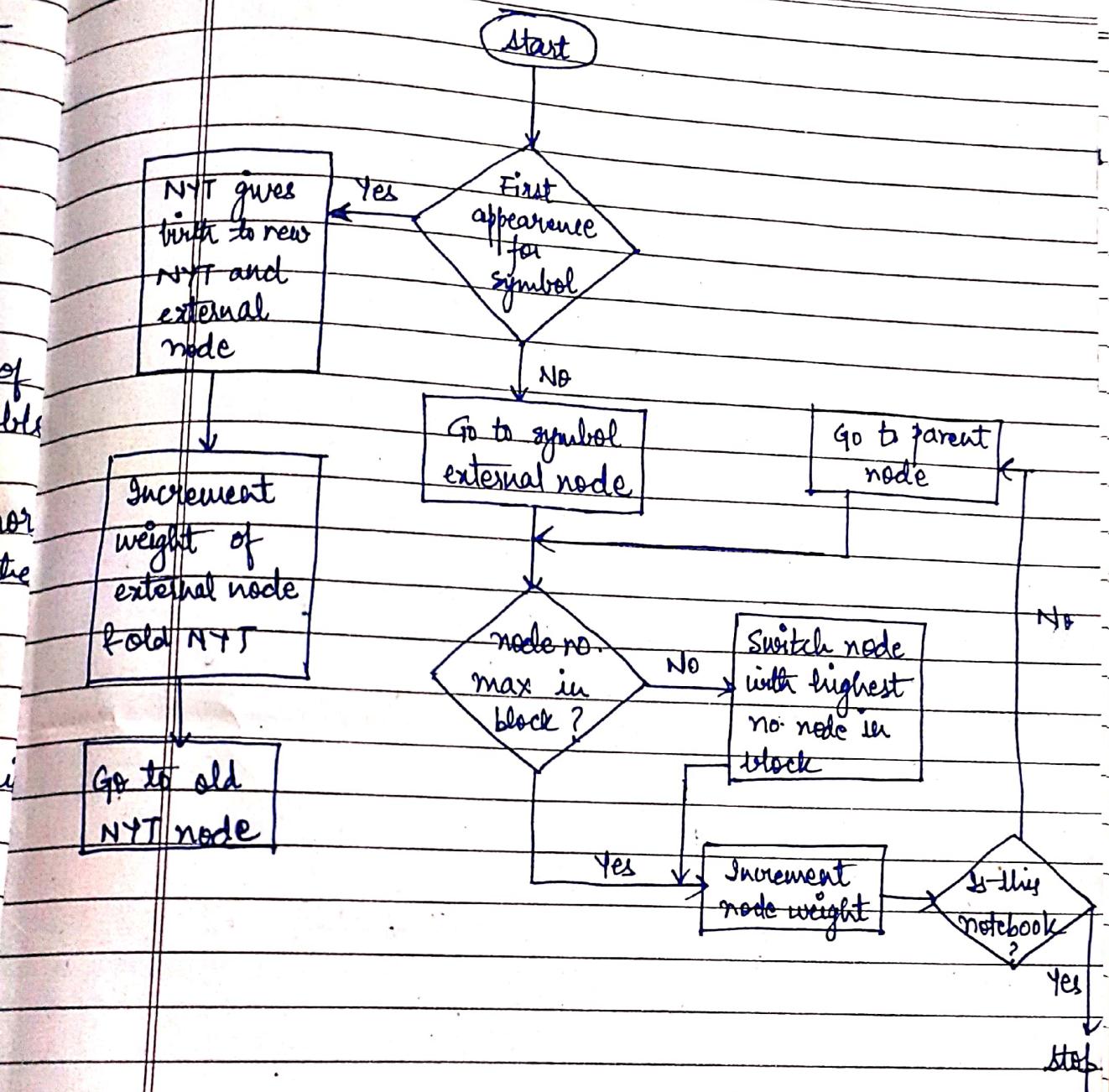
- Update Procedure

Update procedure requires that the nodes be in fixed order. This ordering is preserved by numbering the nodes. Largest code given to root node in the tree and smallest no. is assigned to NYT node.

No. from NYT node to the root are assigned in an increasing order from left to right & from lower to upper layer.



Ex 4



Tree Update (For Encoding for decoding)

Two Parameters : 1. weight of each leaf
2. Node no. (Every node have code no.)

weight of each leaf

For external node

For internal node

weight
no. of time symbol
encountered

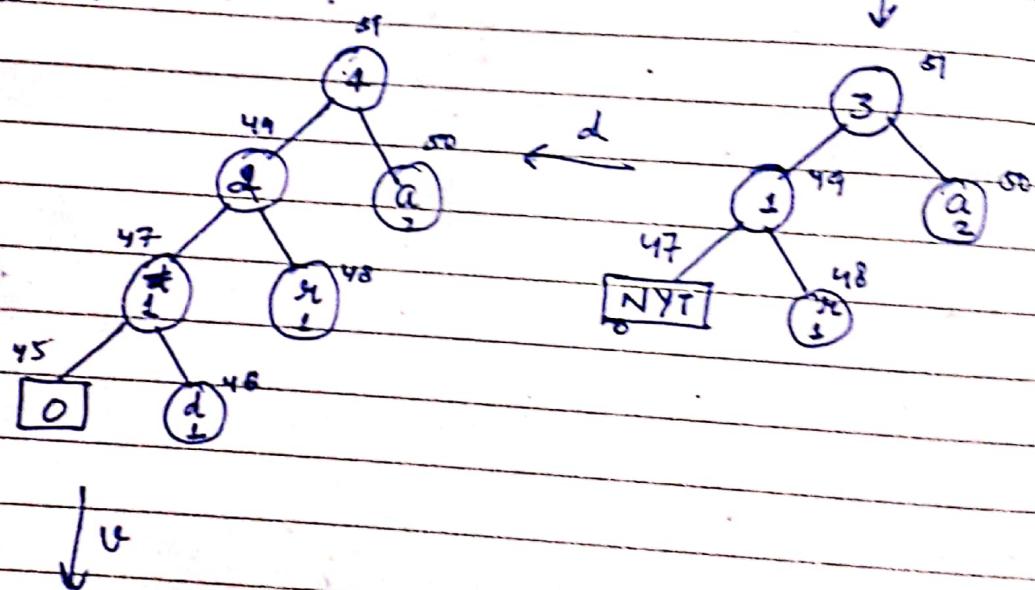
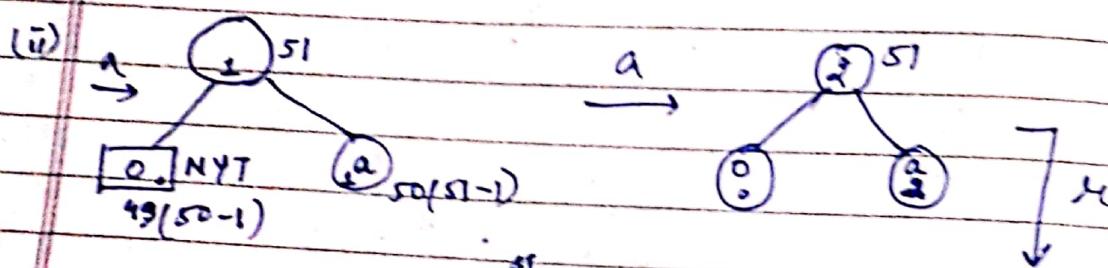
It consist of sum
of weight of its child

node no. = total no. of nodes in tree
 $= (2m-1)$

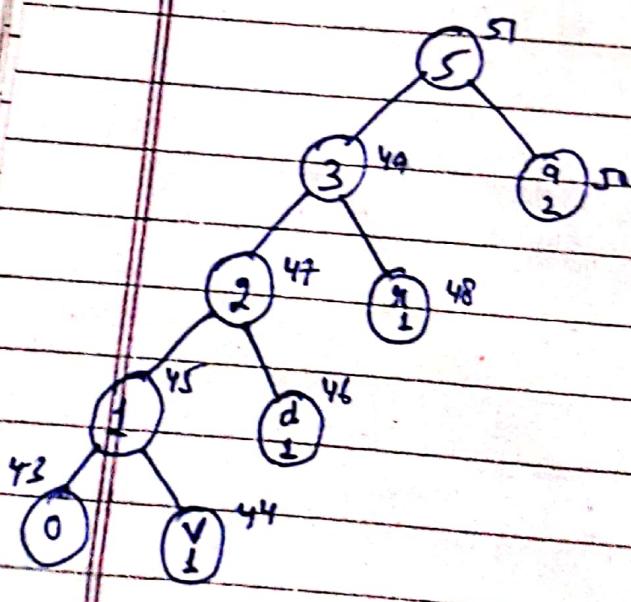
Encoding Procedure

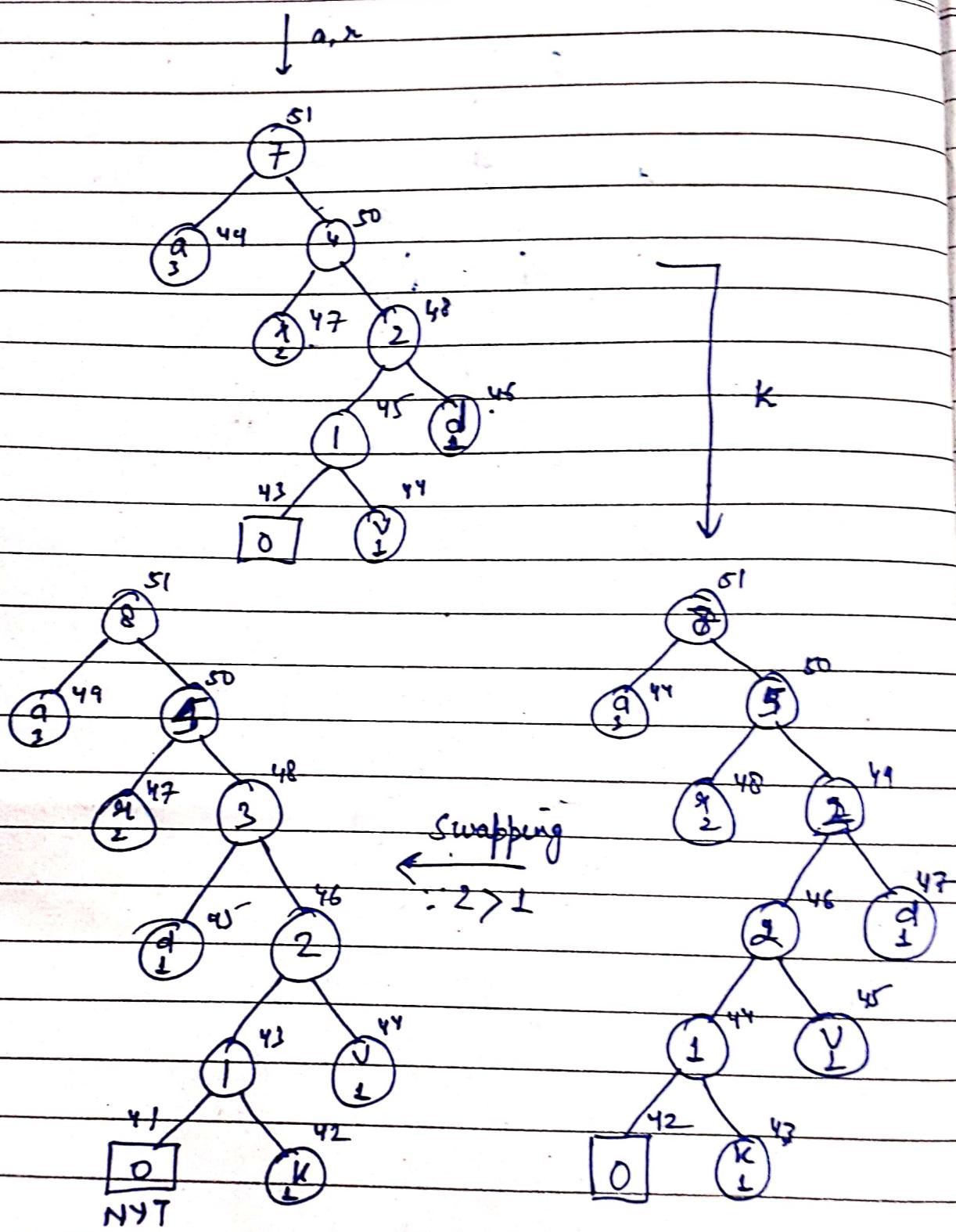
"aardvark" encode with adaptive Huffman coding.

(i) $10 \mid \text{NYT}$ $(2m-1)$ where $m=26$



Swapping when $\text{wt}[\text{left}(\text{child})] < \text{wt}[\text{right}(\text{child})]$





- Adaptive Huffman Coding (Encoding Procedure)

- If the symbol is encountered first time then we encode NYT code followed by fixed.
- If symbol is already encountered, only send code for that symbol.

→ NYT Code

Traversing the tree from root node to NYT node

→ Fixed code

For this we have two parameters e & r

$$m = 2^e + r.$$

m = no. of alphabet in dictionary.

$$m = 26 = 2^e + r.$$

where $e = 4, r = 10$ is fixed.

(i) If ~~$k \neq 1$~~ $1 \leq k \leq 2r$

↑ position of symbol in dictionary
then symbol encoded as $(e+1)$ -bit binary representation of $(k-1)$.

(ii) If $k > 2r$ then symbol encoded as e -bit binary representation of $(k-r-1)$

Q. Encode "a a r d v a n k"

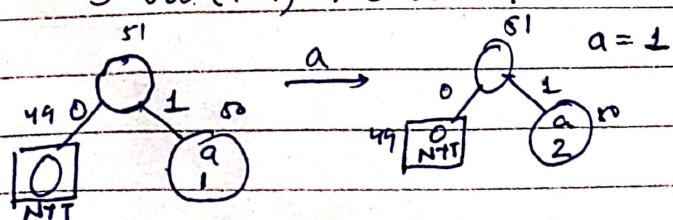
a - NYT - +

$$1 \leq k \leq 2r$$

$$1 \leq 1 \leq 20 \quad (e+1) \text{-bit } (k-1)$$

$$5 \text{ bit. } (1-1) \Rightarrow 5 \text{ bit. } 0.$$

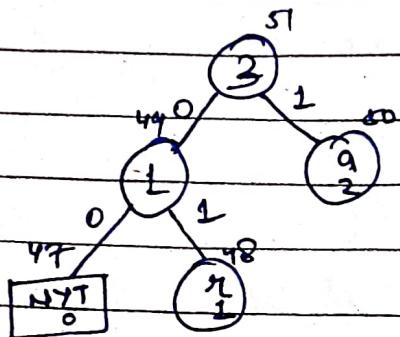
$$a = 00000$$



$$n \rightarrow 1 \leq 18 \leq 20$$

$$r \rightarrow 5 \text{ bit } (18-1) \Rightarrow 5 \text{ bit } 17$$

$$n = 10001 \xrightarrow{\substack{\text{adding} \\ \text{NTT}}} 010001$$

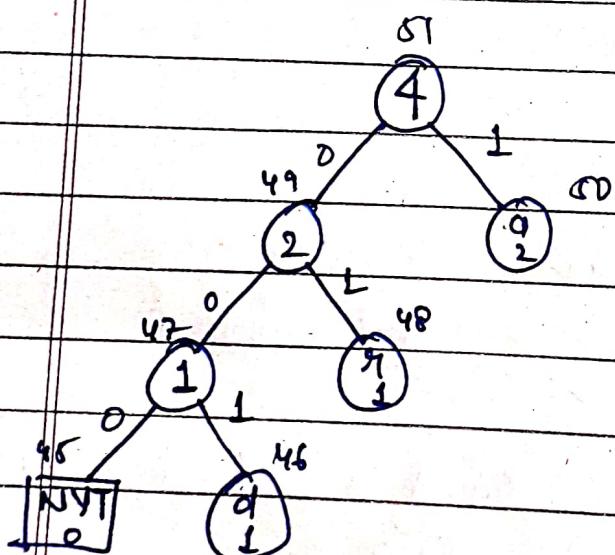


$$d \rightarrow 1 \leq 4 \leq 20$$

$$d \rightarrow 5 \text{ bit } (4-1) \Rightarrow 5 \text{ bit } \#3.$$

$$d = 00011$$

$$d = 0\ 00011$$

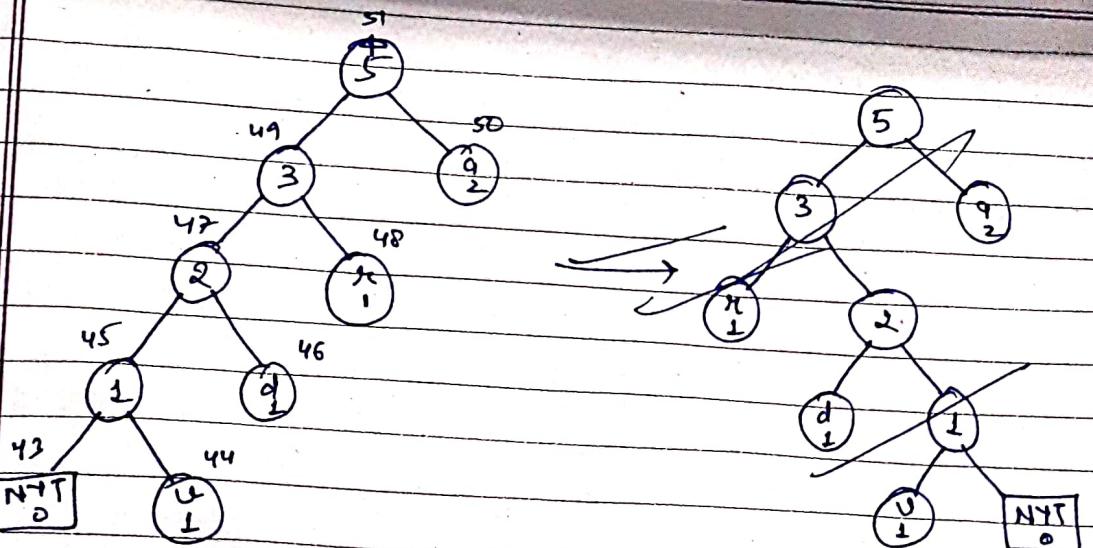


$$v \rightarrow 22 > 20 \text{ (Case ii)}$$

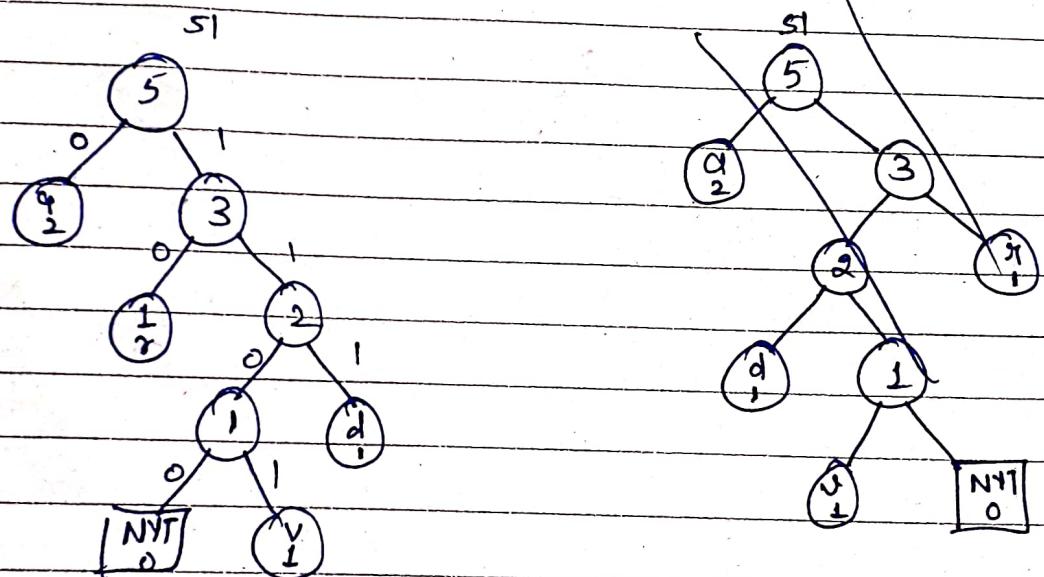
$$v \rightarrow 22 - 10 - 1 = 11 \text{ in 4}$$

$$v \rightarrow 1011$$

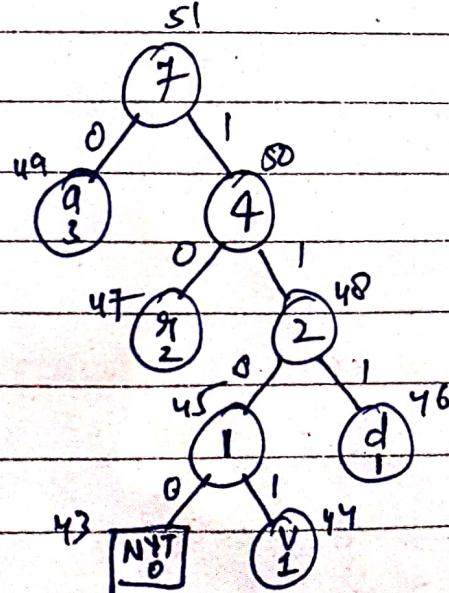
$$NYT \rightarrow 000.$$



swapping



a, g.



$K \rightarrow 51 < 11 < 20$

$K \rightarrow 5 \text{ fits } 10$

$K \rightarrow 01010$

$NYT \rightarrow 1100.$

