

Huffman Codes  
AND  
Data Compression

①

- Computers understand only binary language.
- An encoding scheme maps symbols to binary strings.
- A simple encoding scheme assigns fixed number of bits to each symbol.
- With  $b$  bits you can encode  $2^b$  symbols.
- ASCII is a fixed length encoding scheme.

(2)

Any encoding scheme  
should satisfy following  
two objectives:

- (1) store the data  
    Ccompactly.
  - (2) Reconstruct the  
    data Unambiguously.
-

(3)

## Objective - I (Data Compression)

- Encode frequently occurring letters (eg. e, t, a, o, i, n) with short strings.
- Encode infrequently occurring letters (eg x, z, q) with long strings.

(4)

## Objective - 2 (Unambiguous) reconstruction)

- Use prefix codes.
- A prefix code is one in which the code of one symbol is not a prefix of the code of other symbol.

(5)

- \*  $S$  is the set of symbols.
- \*  ~~$S$~~   $n$  is the # of symbols in the text.
- \*  $\forall x \in S : f_x$  is the fraction of letters equal to  $x$
- \*  $f_{x,n}$  is the # of letters equal to  $x$
- \*  $\sum_{x \in S} f_x = 1$
- \* Let function  $r$  be the P. Code
- \* Encoding Length =  $\sum_{x \in S} n \cdot f_x |r(x)|$
- \*  $ABL(r) = \sum_{x \in S} f_x \cdot |r(x)|$

## EXAMPLE - 1

(6)

$$S = \{a, b, c, d, e\}$$

$$f_a = 0.32, r(a) = 11$$

$$f_b = 0.25, r(b) = 01$$

$$f_c = 0.20, r(c) = 001$$

$$f_d = 0.18, r(d) = 10$$

$$f_e = 0.05, r(e) = 000$$

Calculate the  
average bits per letter (ABL).

$$ABL(r) = 2.25$$

## EXAMPLE - 1

(7)

### IMPROVEMENT IN CODE

$$r(a) = 11$$

$$r(b) = 10$$

$$r(c) = 01$$

$$r(d) = 001$$

$$r(e) = 000$$

The frequency of each letter remains same.

Now calculate ABL(r) again.

(8)

## PROBLEM

- We are given an alphabet  $S$  and a frequency  $f_x$  for each letter  $x \in S$ .
- Produce a prefix code  $\gamma$  for which  $ABL(\gamma)$  is minimum.
- In other words, produce an optimal code  $\gamma$ .

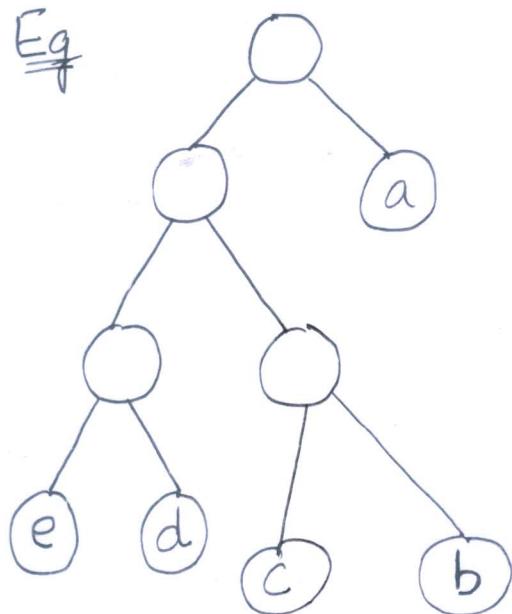
(9)

## PREFIX CODES



## BINARY TREES

Eq



a = 1

b = 011

c = 010

d = 001

e = 000

# of leaves = |S|

(10)

PROBLEM    RESTATE D

We are given an alphabet  $S$ .  
and a frequency for each  
 $x \in S$ . Produce a binary tree  
 $T$  s.t.  $\sum_{x \in S} f_x \cdot \text{depth}(x)$  is MIN.

OR.

PRODUCE an optimal Tree.

$$ABL(T) = \sum_{x \in S} f_x \cdot \text{depth}(x)$$

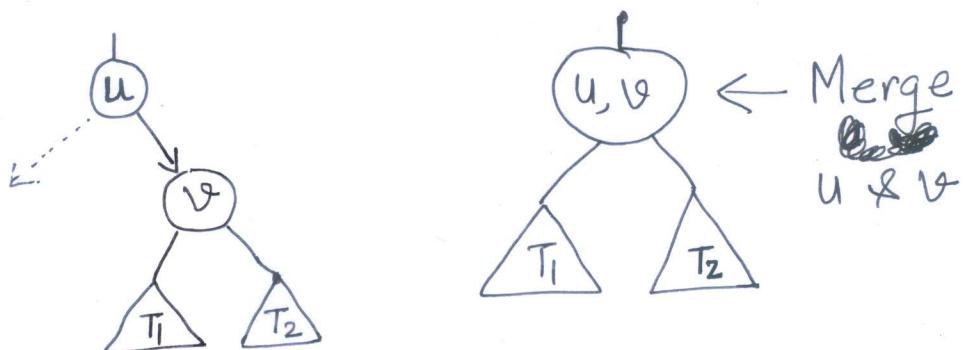
(11)

## LEMMA - I

THE OPTIMAL BINARY TREE

$T^*$  IS A PROPER BINARY TREE.

Proof: If not, you can convert it into one whose ABL value will be less than  $T^*$ , a contradiction. !



(12)

## LEMMA - 2

SUPPOSE  $y$  and  $z$  are leaves of an optimal binary tree.

If  $\text{depth}(y) < \text{depth}(z)$ ,

then  $f_y \geq f_z$

Proof : By Contradiction

(13)

## THEOREM

IN OPTIMAL BINARY TREE

$T^*$ , THE TWO LOWEST  
FREQUENCY LETTERS ARE  
ASSIGNED TO LEAVES  
THAT ARE SIBLINGS.

USING THIS THEOREM THE  
TREE CAN BE CONSTRUCTED  
AS FOLLOWS:

(14)

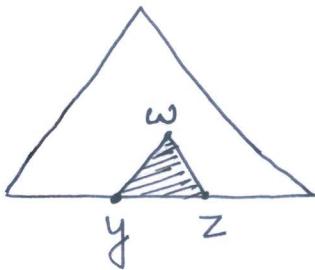
- Initially the Algo is called on alphabet size  $k$
- In the second step, the Algo invokes itself on alphabet size  $k-1$
- In the final step, the Algo invokes itself on alphabet size 2.

## THE FIRST STEP

(15)

① In the 1st step, the two lowest frequency letters  $y$  &  $z$  are merged into a single letter  $w$ .

② The size of the alphabet reduces from  $k$  to  $k-1$ .



③  $T$  and  $T'$  are Trees built from alphabet size  $k$  and  $k-1$ , resp.

- $f_w = f_y + f_z$
- $\text{depth}_{T'}(x) = \text{depth}_T(x) \begin{cases} x \neq y \wedge \\ x \neq z \end{cases}$
- $\text{depth}_{T'}(w) = \text{depth}_T(y) - 1$
- $\text{depth}_{T'}(w) = \text{depth}_T(z) - 1$

(16)

Theorem  $ABL(T') = ABL(T) - f_w$

Proof:

$$ABL(T) = \sum_{x \in S} f_x \cdot \text{depth}_T(x)$$

$$= f_y \cdot \text{depth}_T(y) + f_z \cdot \text{depth}_T(z)$$

$$+ \sum_{x \neq y, z} f_x \cdot \text{depth}_T(x)$$

$$= (f_y + f_z) (1 + \text{depth}_{T'}(\omega))$$

$$+ \sum_{x \neq y, z} f_x \cdot \text{depth}_{T'}(x)$$

$$= f_w + f_w \cdot \text{depth}_{T'}(\omega)$$

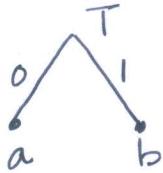
$$+ \sum_{x \neq y, z} f_x \cdot \text{depth}_{T'}(x)$$

$$= f_w + ABL(T')$$

## PROOF OF CORRECTNESS

(17)

- We prove by MI on size of the alphabet
- **Base Case**      Size = 2       $S = \{a, b\}$  say



T is optimal

- **Induction Step**      Suppose the tree built when size = K is sub optimal.

- It means,  $\exists$  tree  $Z$  s.t  $ABL(Z) < ABL(T)$
- Convert  $Z$  into  $Z'$  by removing two least frequency symbols  $y$  and  $z$ , and labeling their parent as  $w$
- Note both  $y$  and  $z$  were leaves with common parent.
- Now by previous theorem  
 $ABL(Z) = ABL(Z') + fw$   
 $ABL(T) = ABL(T') + fw$

(18)

- Since  $A\overline{B}L(Z) < A\overline{B}L(T)$ ,  
we have  $A\overline{B}L(Z') < A\overline{B}L(T')$ .
- $T'$  is suboptimal for  
alphabet of size  $k-1$ .
- What we have proved is :  
If  $T$  is suboptimal for alphabet  
of size  $k$ , then  $T'$  is  
suboptimal for alphabet of  
size  $k-1$ .
- The contrapositive of above is :  
If  $T'$  is optimal for alphabet of  
size  $k-1$ , then  $T$  is optimal  
for alphabet of size  $k$ .
- The Induction Step is over.

(19)

## COMPLEXITY OF THE HUFFMAN ALGO

- (1) Construct a heap of the  $k$  letters based on their frequencies :  $O(k)$
- (2) Fetch two min freq letters  
Say  $x$  and  $y$  from  
the heap:  $O(\log k)$ .
- (3) Add a new letter  $z$  into the  
heap s.t  $f_z = f_x + f_y$  :  
 $O(\log k)$
- (4) Repeat 2, 3  $k$ -times.  
 $O(k \cdot \log k)$