Class Partipation — No credit.

Hufmann Encoding
- Why we study it.
- Formal def$^n$
- Greedy Algorithm

**Encoding:** Representing symbols (in message) using binary codes 0/1.

Eg. msg = C A B D C

| SYMBOL | CODES |
|--------|-------|
| A | 00 |
| B | 01 |
| C | 10 |
| D | 11 |

2 bit string

Eg. In computer code (ASCII co...

A — 01 00001  (65)
B — 01 00010  (66)
⋮

8 bit string

**"Variable Length Coding"** — All codes need not be of same length.

Eg.

Table →

| SYMBOLS | CODES |
|---------|-------|
| A | 0 |
| B | 100 |
| C | 101 |
| D | 11 |

length ∈ [1, 3]

Question — When can such encodig be useful.

Eg. freq: A−45, B−9, C−11, D−35 in a msg of 100 char.

100 = 45 + 9 + 11 + 35

Length of encoded msg. = 45(1) + 9(3) + 11(3) + 35(2)

= 175 ← 12% impro...

If we were using 00, 01, 10, 11 as encodig, then length = 200

len = 2

Is there ambiguity?

C
⇑
... → 10 1 ... ... ← ...

msg. = CABAD → 101 0 100 0 11 ←


Not mapped
Not mapped
mapped to C

To take care of ambiguity we ensure

if $(x_1 \ldots x_k)$ is CODE-WORD, then no prefix of it is co[

Prefin Encoding : ⟵ defⁿ

└ Eg. Country call nos.

$+91$ — India
$\frac{+91}{X}$  $+1$ — USA. (10 digit)

$+19$ — Not a country code. ⟩ Ambiguity
(9 digit)

**Huffman Encod-**
ⓐ JPEG / MP3
   compressi
ⓑ Zipping a file

## Tree Representation

| | CODES |
|---|---|
| A | 0 |
| B | 100 |
| C | 101 |
| D | 11 |



Property of Prefix
→ All symbols m
   be leaf nodes

← length = 4 → 3

## Problem :

Given : Symbols $(a_1 \ldots a_n)$ with freq. vector $F = (f_1 \ldots$

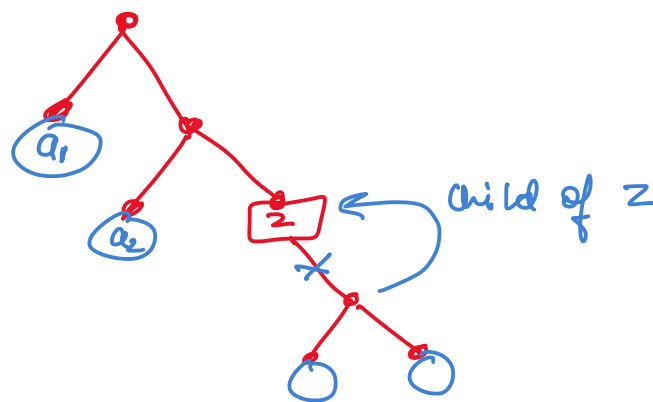Find : Prefix encoding for which encoded msg has "min"

A binary tree T with leaves $(a_1 \ldots a_n)$ such that

$$\left( \sum_{i=1}^{2} (f_i)(\text{depth of } a_i \text{ in } T) \right) \text{ be } m$$

**Property 1:** Each internal node should have 2 children

**Proof:** By contradiction

Take a tree $T$ which is not complete & let internal node of degree 1.



child of $z$

All leaf
$\longrightarrow$ still co

$\sum_{i=1} f_i$
is re

**Ques:** Is it necessary that one child of each no

No.
A — 25
B — 25
C — 25
D — 25



**Property 2:** If $f_1 \geqslant f_2 \geqslant \cdots \geqslant f_n$. Then in o

① $\text{depth}(a_1) \leq \text{depth}(a_2) \leq \cdots \leq \underline{\text{depth}(a_n)}$

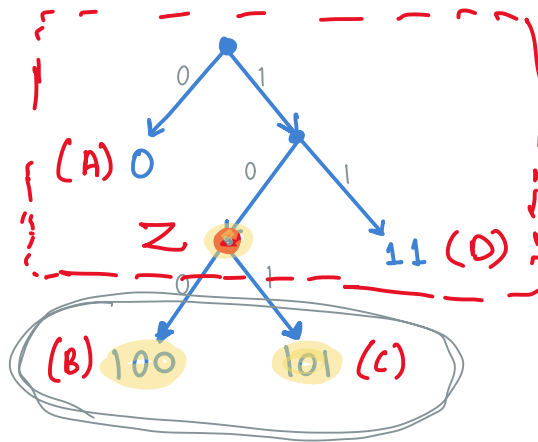② $\text{depth}(a_n) = \text{depth}(a_{n-1})$

**Proof of ②:** Both children of parent (an)

MAX —

are leaf node, & will correspond to _two_ minimum frequencies | output

**CODES**

| | |
|---|---|
| A | 0 |
| ~~B~~ | 100 |
| ~~C~~ | 101 |
| D | 11 |

Z



(A) 0

Z

11 (O)

(B) 100    101 (C)

(F)

A — 45
B — 9 }
C — 11 }
D — 35
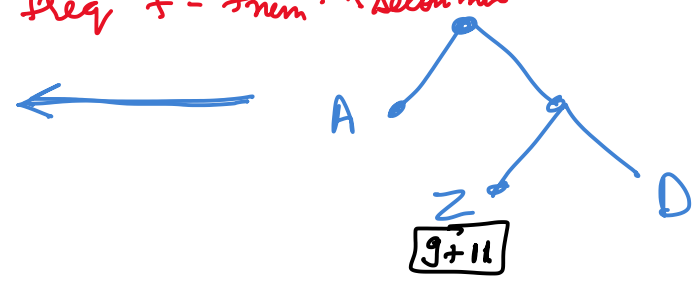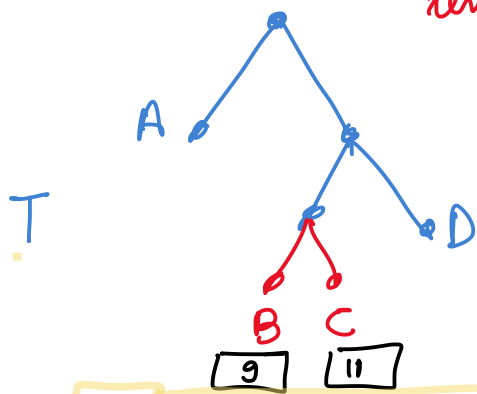
Greedy strategy →

Take minimum
2 freq & replace
term with new
freq $f = f_{min} + f_{second\ min}$

F'

A — 45
Z — 20 }
D — 35

↓ tree T'

T

A

B   C
[9] [11]

D

← 

A

Z
[9+11]

D

$$|encoded-msg(F)| = |encoded-msg(F')| + (1) ($$

diff in le

**H.W.** — Correctness of this algo.

**Theorem:** If $a_{n-1}, a_n$ have least frequency, then

problem

$F^* = (F \cup \{\ldots\}) \setminus \{\ldots$

$$\Gamma = (F \cup \{ t_n + t_{n-1} \}) \setminus \{ t_n, t_n \}$$

$$opt(F) = opt(F^*) + (f_{n-1} + f_n).$$

H.W. — Find an $O(n \log n)$ time implementation.