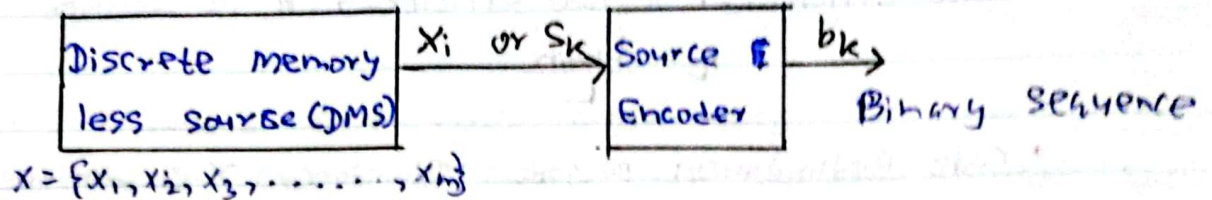Communication Theory II
Lecture 07
Source Coding Theorem, Huffman Coding

Source Coding



$$X_i \text{ or } S_k$$

Discrete memory less source (DMS) → Source Encoder → $b_k$ Binary Sequence

$$X = \{x_1, x_2, x_3, \ldots \ldots, x_m\}$$

• A conversion of the output of a discrete memoryless source (DMS) into a sequence of binary symbols (i.e., binary code word) is called Source Coding.

• The device that performs this conversion is called the source encoder. Above figure shows a source encoder.

• There are few terms related to Source Coding Process:
   i. Code word length
   ii. Average Code word length
   iii. code efficiency
   iv. Code redundancy

• Code Word Length: Let x be a DMS with finite entropy $H(X)$ and an alphabet $\{x_1, \ldots x_m\}$ with corresponding probabilities of occurrence $P(x_i)$ $(i = 1, \ldots, m)$.

• Let the binary code word assigned to symbol $x_i$ by the encoder have length $n_i$, measured in bits

• The length of a code word is the number of binary digits in the codeword.

• **Average Code Word Length:** The average Code word length per source is given by

$$L = \sum_{i=1}^{m} P(x_i) n_i$$

• **Code Efficiency:** The code efficiency $\eta$ is defined as under:

$$\eta = \frac{L_{min}}{L}$$

• **Code Redundancy:** The code redundancy $\gamma$ is defined as

$$\gamma = 1 - \eta$$

## Source Coding Theorem

• The Source coding theorem states that for a DMS $X$, with entropy $H(X)$, the average code word length $L$ per symbol is bounded as $L \geq H(X)$

• Further, $L$ can be made as close to $H(X)$ as desired for Some suitably chosen one.

• Thus, with $L_{min} = H(X)$ the code efficiency can be written as
$$\eta = \frac{H(X)}{L}$$

• **Classification of code:** Classification of codes is best illustrate by an example.
   1. Fixed-Length codes
   2. Variable-Length codes
   3. Distinct codes
   4. Prefix-Free codes
   5. Uniquely Decodable codes
   6. Instantaneous codes
   7. Optimal codes

## Table. Binary Codes

| $X_i$ | Code 1 | Code 2 | Code 3 | Code 4 | Code 5 | Code 6 |
|-------|--------|--------|--------|--------|--------|--------|
| $X_1$ | 00 | 00 | 0 | 0 | 0 | 1 |
| $X_2$ | 01 | 01 | 1 | 10 | 01 | 01 |
| $X_3$ | 00 | 10 | 00 | 110 | 011 | 001 |
| $X_4$ | 11 | 11 | 11 | 111 | 0111 | 0001 |

- Code 1 & Code 2 are fixed-length codes
- Code 3, Code 4, Code 5 & Code 6 are variable-length codes
- If code is distinct, each codeword is distinguishable from other codewords
- If code inputs no codewords can be found by adding code symbols to another codeword is called Prefix-free code
- In prefix-free code, no codewords are is prefix of another codeword
- Code 2, Code 4 & Code 6 are Prefix-Free codes

- Code 3 of the table is not uniquely decodable. For an example, in binary code sequence. 1001 may correspond to the source sequence $X_2 X_3 X_2$ or $X_2 X_1 X_1 X_2$

- Uniquely decodable codes are called as instantaneous codes.
- The instantaneous codes have the property of previously mentioned that no codeword is prefix of another codeword.
- For this reason, Prefix free codes are sometimes known as instantaneous codes.

- If code is set to be optimal, if it is instantaneous and has minimum average for a given source, with a given probability for the source symbol.
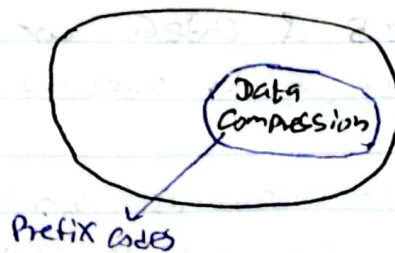
# Prefix-Free Codes

A Prefix-free Code is one in which no Codeword is a prefix of any other. Example

### Explain Prefix Codes

Where do Prefix Codes occur?

Coding Theory



Prefix Codes

## Why Prefix Codes are important

Because

• Unique decodability → This a must have
→ Inst

Imagine

| letters | Code |
|---------|------|
| a | 0 |
| b | 1 |
| c | 00 |

decode: 100

It is 'baa' or 'cc'...

Instantaneous decoding → which is a nice-to-have

Imagine

| letters | Code |
|---------|------|
| a | 0 |
| b | 01 |
| c | 011 |

decode: 011

1 option: 'c'
we need to keep decoding (reading)
until another zero occurs

## Definition

- A prefix code is a code in which no codeword is a prefix (initial segment) to another codeword

| Letters | Code |
|---------|------|
| a | 0 |
| b | 10 |
| c | 11 |

- Prefix codes are also known as:
  - Prefix-free Codes
  - Prefix Condition Codes
  - Instantaneous codes
  - Huffman codes

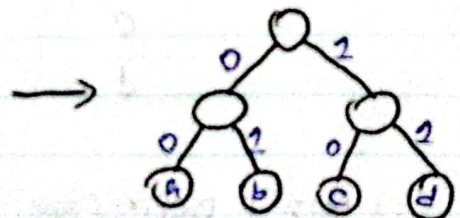| Letters | Code |
|---------|------|
| a | 00 |
| b | 01 |
| c | 10 |
| d | 11 |

- Fixed-length codes are always prefix-codes

- To check for any code if it is a prefix code, draw its binary tree and check if all codewords are external nodes (i.e. leaves)

| Letters | Code |
|---------|------|
| a | 0 |
| b | 10 |
| c | 11 |



## How to recognize Prefix codes

| Letters | Code |
|---------|------|
| a | 00 |
| b | 01 |
| c | 10 |
| d | 11 |

$\longrightarrow$



all are external nodes

Is this a Prefix Code?

Yes, since this is a fixed-length code.

| Letters | Code |
|---------|------|
| a | 1 |
| b | 00 |
| c | 01 |
| d | 10 |
| e | 11 |



internal node

## 2. Is this a prefix code?

No, 1 is a prefix to 10 and 11.
(no unique decodability)

| Letters | Code |
|---------|------|
| a | 0 |
| b | 01 |
| c | 011 |
| d | 0111 |



internal node
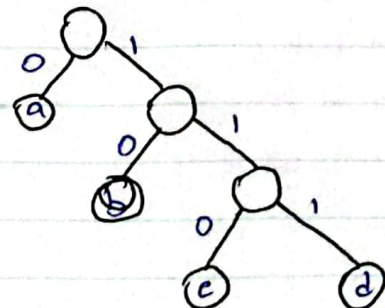
## 3. Is this a prefix code?

No, 0 is a prefix to 01, 011 and 0111,
01 is a prefix to 011 and 0111, and 011
is a prefix to 0111
(no instantaneous decoding)

| Letters | Code |
|---------|------|
| a | 0 |
| b | 10 |
| c | 110 |
| d | 1110 |



Is this a prefix code?

Yes!

all are external nodes

# Prefix-Free codes [Example 2]

A Prefix-Free code is one in which no codeword is a prefix of any other

An Example

| | |
|---|---|
| Space | 1 |
| E | 01 |
| T | 001 |
| A | 0001 |
| O | 00001 |
| I | 000001 |

I ATE: 000001100100101





| | |
|---|---|
| Space | 111 |
| E | 010 |
| T | 1101 |
| A | 1011 |
| O | 1001 |
| I | 1000 |

I ATE: 100011110111101010

$$\text{Information} \quad = \quad \log_b \frac{1}{P(n_i)} \quad = \quad -\log_b P(n_i)$$
$$I(x_i)$$

$$\text{Entropy} \quad = \quad \sum_{k=1}^{h} P_k \log_2 \frac{1}{P_k} \quad = \quad -\sum P_k \log_2 P_k$$

Average code word length
$$L = \sum_{k=1}^{b} P_k n_k$$

$P \Rightarrow$ Probability $\qquad\qquad n = $ number of bits

| m | P | Stage 1 | Stage 2 | Stage 3 | codeword | length |
|---|-----|---------|---------|---------|----------|--------|
| $m_1$ | 0.30 | 0 | 0 | | 0 0 | 2 |
| $m_2$ | 0.25 | 0 | 1 | | 0 1 | 2 |
| $m_3$ | 0.15 | 1 | 0 | 0 | 1 0 0 | 3 |
| $m_4$ | 0.12 | 1 | 0 | 1 | 1 0 1 | 3 |
| $m_5$ | 0.10 | 1 | 1 | 0 | 1 1 0 | 3 |
| $m_6$ | 0.08 | 1 | 1 | 1 | 1 1 1 | 3 |

entropy $\Rightarrow$ $H = \sum_{k=1}^{6} P_k \log_2 \frac{1}{P_k}$

$= (0.30) \log_2 \frac{1}{0.30} + (0.25) \log_2 \frac{1}{0.25} + (0.15) \log_2 \frac{1}{(0.15)}$

$+ (0.12) \log_2 \frac{1}{(0.12)} + (0.10) \log_2 \frac{1}{(0.10)} + (0.08) \log_2 \frac{1}{0.08}$

$= 2.418$ bits

Average codeword length :
$$L = \sum_{k=1}^{6} P_k n_k$$

$= (0.30 \times 2) + (0.25 \times 2) + (0.15 \times 3) + (0.12 \times 3) + (0.10 \times 3)$
$\quad + (0.08 \times 3)$

$= 2.45$ bits

## Efficiency

$$\eta = \frac{L_{min}}{L}$$

If $\eta = 1$ then Source code is considered as

**efficient**

$$\eta = \frac{H}{L}$$

## Huffman Coding Algorithm

~~Example: Alphabet~~

**Steps:**

1. The Source symbols are arranged in order of decreasing probability. Then the two of lowest Probability are assigned bit 0 and 1.

2. Then combine last two symbols and move the combined symbol as high as possible.

3. Repeat the above step until end.

4. Code for each symbol is found by moving backward.

5. eta calculation.

**Efficiency**

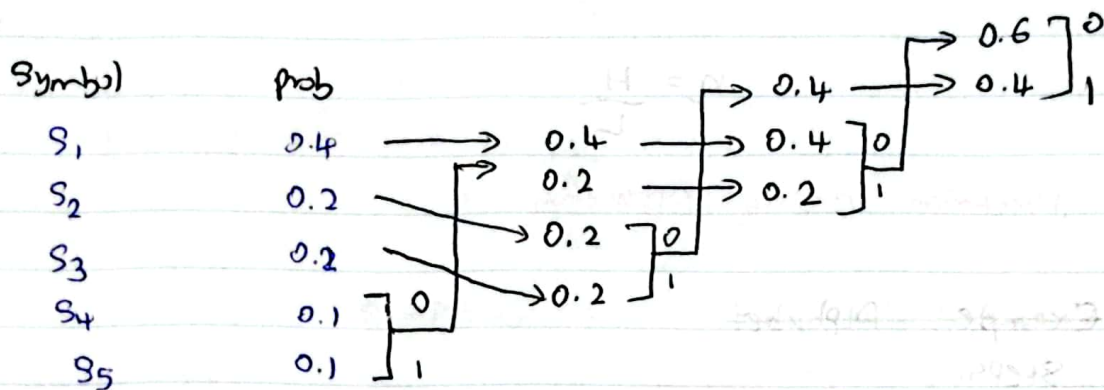$$\eta = \frac{H}{L \log_2 r} \longrightarrow \text{binary } r = 2 \ (0,1)$$

$$L = Avg \ Codeword = \sum_{i=1}^{n} P_i \ n_i$$

$$H = \sum_{i=1}^{n} P_i \log_2 \left(\frac{1}{P_i}\right)$$

$$Variance \quad \sigma^2 = \sum_{i=1}^{n} P_i \ (n_i - L)^2$$

**Example:**

Alphabet with prob = {0.4, 0.2, 0.2, 0.1, 0.1} For Symbols {$S_1$, $S_2$, ..., $S_5$}. Find Huffman Codes and also find efficiency & variance.

| Symbol | prob |
|--------|------|
| $S_1$ | 0.4 |
| $S_2$ | 0.2 |
| $S_3$ | 0.2 |
| $S_4$ | 0.1 |
| $S_5$ | 0.1 |

$$0.4 \rightarrow 0.4 \rightarrow 0.4 \rightarrow 0.6 \rbrack^0$$
$$0.2 \rightarrow 0.2 \rightarrow 0.4 \rbrack^0 \rightarrow 0.4 \rbrack_1$$
$$0.2 \rightarrow 0.2 \rbrack^0 \rightarrow 0.2 \rbrack_1$$
$$0.1 \rbrack^0 \rightarrow 0.2 \rbrack_1$$
$$0.1 \rbrack_1$$

| Symbol | Codeword | length |
|--------|----------|--------|
| $S_1$ | 00 | 2 |
| $S_2$ | 10 | 2 |
| $S_3$ | 11 | 2 |
| $S_4$ | 010 | 3 |
| $S_5$ | 011 | 3 |

**—Entropy**

$$H = \sum P_i \log_2\left(\frac{1}{P_i}\right)$$
$$= 0.4 \log_2\left(\frac{1}{0.4}\right) + 2 \times 0.2 \log_2\left(\frac{1}{0.2}\right) + 2 \times 0.1 \log_2\left(\frac{1}{0.1}\right)$$
$$= 2.1216 \quad bits/symbol$$

**— $L = \sum P_i n_i$**

$$= 2 \times 0.4 + 2(2 \times 0.2) + 3 \times 0.1 \times 2$$
$$= 2.2 \quad bits/symbol$$

$$\eta = \frac{H}{L \log_2 r} = \frac{2.1216}{2.2 \times \log_2 2} = 96.4\%$$

$$\sigma^2 = \sum P_i (n_i - L)^2$$
$$= 0.4(2 - 2.2)^2 + 2 \times 0.2(2 - 2.2)^2 + 2 \times 0.1(3 - 2.2)^2$$
$$= 0.16 \quad \text{As Low as possible}$$

# Huffman Algorithm

Two types of encoding

→ fixed length encoding
  ↳ assigns to each character a bit string
    of the same length m

→ Variable length encoding
  ↳ assigns codewords of different lengths to
    different characters
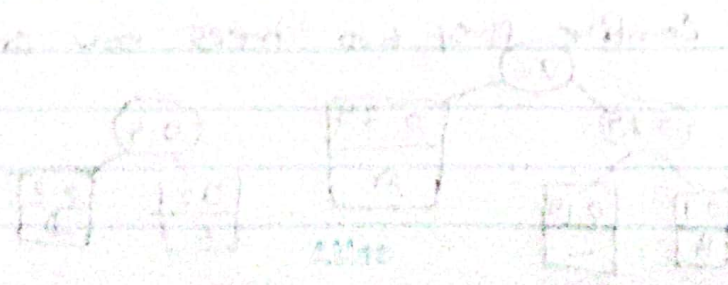    → Huffman algorithm coming under
      variable length encoding

## Huffman Algorithm

1. Initialize n one node trees and label them
   with the characters of the alphabet.

2. Record the frequency of each character in it's
trees root to indicate trees weight

3. Repeat the following operation until a single
tree is obtained. Find two trees with the smallest weight.
make them the left and right subtree of a new tree and
record the sum of their weights in the root of the new tree
as its weight

   — A tree constructed by the above algorithm is called
a Huffman tree.

   Ex

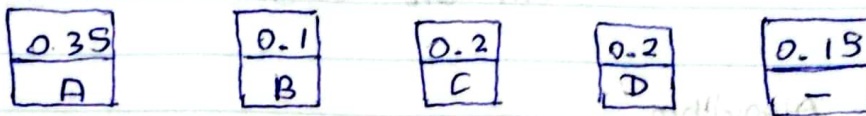   consider the five character alphabet A, B, C, D, _ with
the

ex

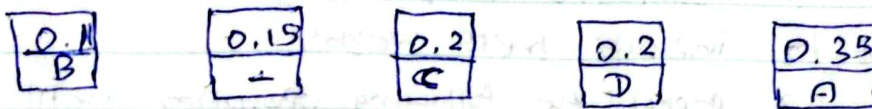Consider the five character alphabet $A, B, C, D, \_$ with the following occurrence Probability.

| Character | A | B | C | D | \_ |
|-----------|------|-----|-----|-----|------|
| Probability | 0.35 | 0.1 | 0.2 | 0.2 | 0.15 |

## Step 1

Create single node tree

| 0.35 | | 0.1 | | 0.2 | | 0.2 | | 0.15 |
|------|---|-----|---|-----|---|-----|---|------|
| A | | B | | C | | D | | \_ |

arrange the single node trees with ascending order

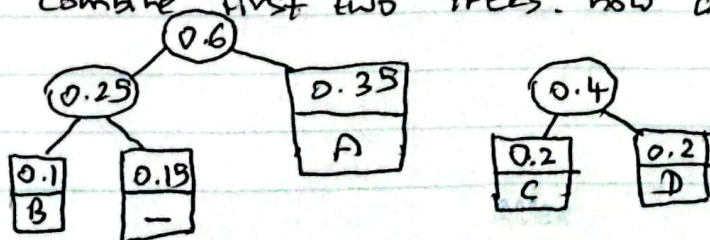| 0.1 | | 0.15 | | 0.2 | | 0.2 | | 0.35 |
|-----|---|------|---|-----|---|-----|---|------|
| B | | \_ | | C | | D | | A |

combine first two trees and create new tree
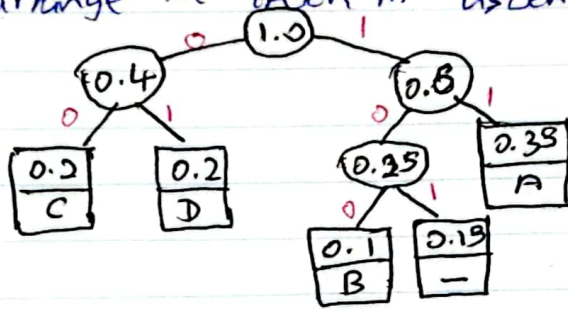


again arrange trees in ascending order



Combine first two trees. now arrange by ascending order

→ Combine first two trees. arrange

→ arrange the tree in ascending order



→ now mark as left tree as 0 and right tree as 1

→ how "DAD" is encoded as

      011101

→ In the decoding side

     DAD