# CODING FOR DISCRETE SOURCES

Three major types of data sources:

a) analog : continuos wave

b) analog sequence: continuos in value
discrete samples

c) discrete : symbols from finite alphabet $X$

Ex: $X = \{A, B, C, D, \ldots, Z, 0, 1, \ldots, 9, \ldots\}$

$Y = \{ \cdot, \mid, -, \vdash \}$

$X = \{0x00 \ldots 0xff\}$

## Coding

- Source alphabet $A$
- code alphabet $B$

- coding: symbols from $A \longrightarrow$ words in $B$

- word: nonempty finite sequence of symbols

Ex: words $00000, 00001, \ldots, 11111$
10 of them contain two symbols 1
$\Rightarrow$ 2-out-of-5 code

$0 \longrightarrow 00011$

$1 \longrightarrow 11000$

$2 \longrightarrow 10100$

3

$\vdots$

$211 \longrightarrow 10100, 11000, 11000$

### Unique decodability

a) $a_1 \neq a_2 \in A \longrightarrow K(a_1) \neq K(a_2) \in B$

b) no two combinations of source symbols produce the same sequence of code words

Ex: $A = \{0, 1, 2, \ldots, A, B, \ldots, F\}$

$B = \{0, 1\}$

$K(A): \quad 0 \longrightarrow 0000$

$\qquad 1 \longrightarrow 0001$

$\vdots$

$\qquad A \longrightarrow 1010$

$\vdots$

$\qquad F \longrightarrow 1111$

## Fixed-length codes

block codes
→ length of all code words is const.

$A = \{a_1, a_2, \dots, a_M\}$ .... $M$ symbols

$B = \{0, 1\}$

(?) length of a code word $M$

$$M = \lceil \log_2 M \rceil \qquad M \le 2^M$$

## Variable-length codes

Ex: $A = \{a, b, c\}$

$K(a) = 0$
$K(b) = 10$
$K(c) = 11$

Ex: Morse

E → •
T → —
A → • —
Y → — • — —

– codewords of different lengths
– buffering

→ Prefix-free codes

a) source codeword can be decoded as soon as its last bit arrives
b) if a uniquely-decodable code with the given code lengths exists it can be made prefix
c) given the probability distribution of the source symbols, it is possible to design
   a prefix-free code with optimum code-word lengths

Ex: $A \{a, b, c\}$

$a \to 0 \qquad babc \to 10110$
$b \to 1$
$c \to 10$

babc | cb

babba

Ex:

$a \to 1 \qquad a \to 1$
$b \to 10 \qquad b \to 01$
$c \to 100 \qquad c \to 001$

# Construction of a prefix-free code

$\mathcal{A} = \{a, b, c, d, e, f\}$

$a \to 000$
$b \to 001$
$c \to 010$
$d \to 011$
$e \to 10$
$f \to 11$



$a \to 00$
$b \to 010$
$c \to 0110$
$d \to 0111$
$e \to 10$
$f \to 11$

$\longrightarrow$ all source symbols form a leaf of a binary tree
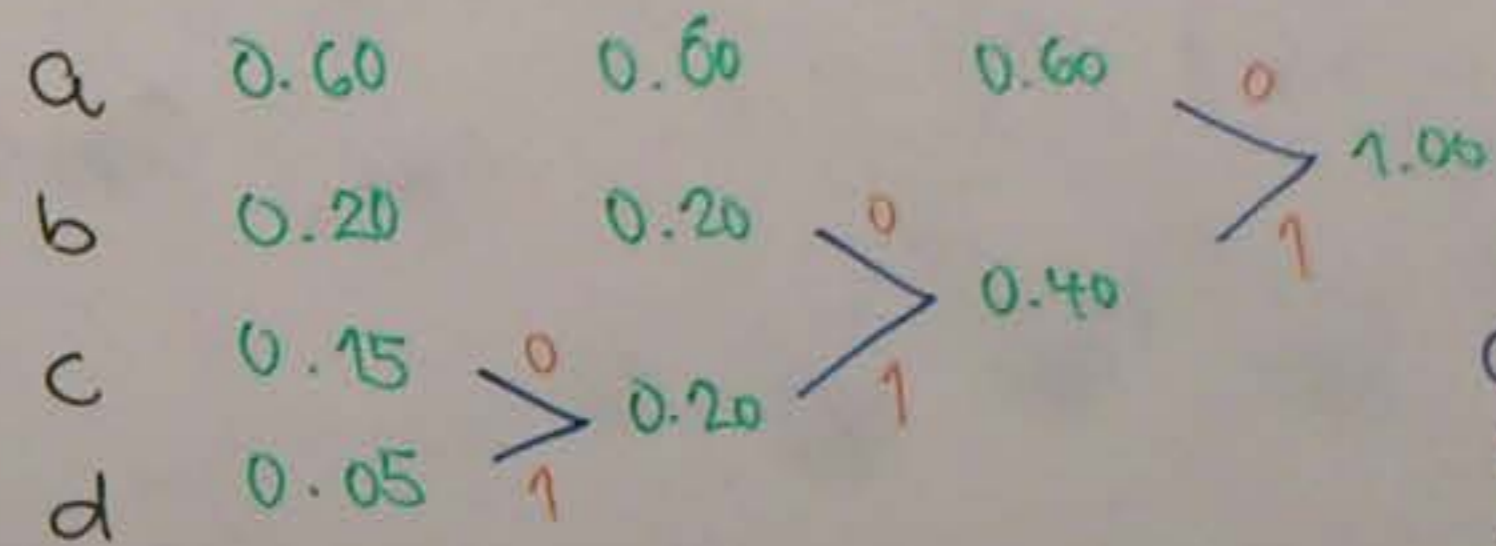
$\longrightarrow$ travelling the the tree from root produces the code word
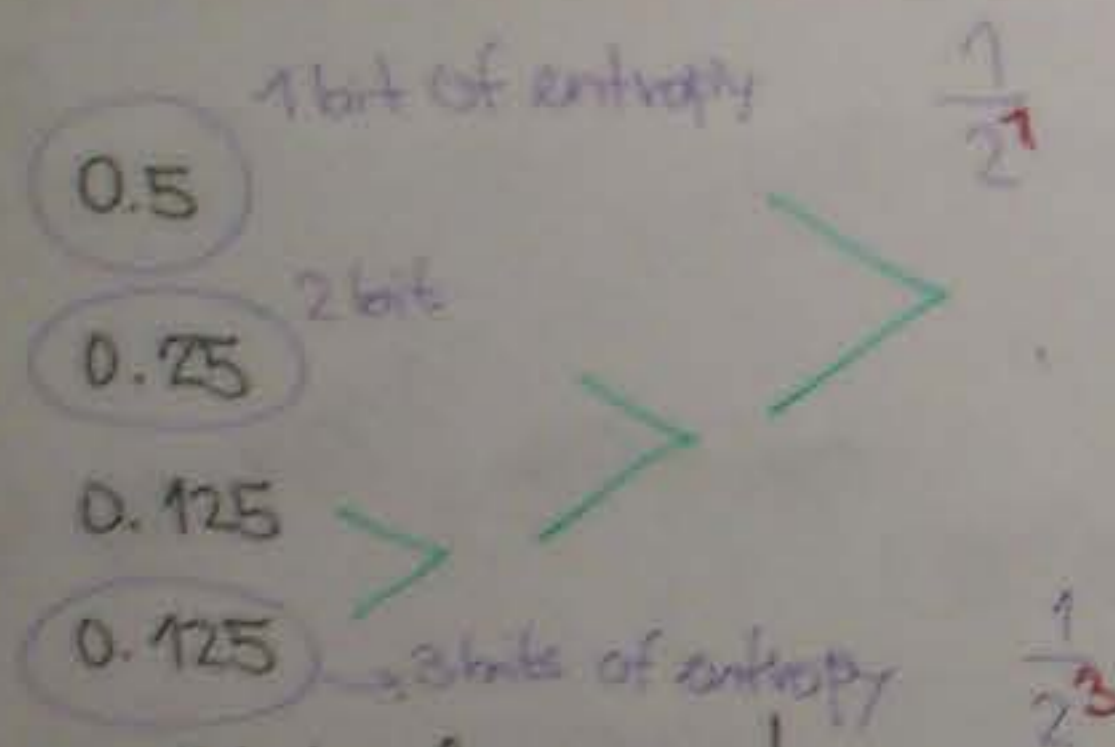
## Optimal code word length

$\mathcal{A} = \{a_1, a_2, \ldots, a_n\}$ + probabilities

$p(a_1), p(a_2), \ldots, p(a_n)$

$\to$ code word lengths approximating $-\log_2 p(a_i)$

## Huffman's algorithm

| | | | | |
|---|---|---|---|---|
| a | 0.60 | 0.60 | 0.60 | 0 |
| b | 0.20 | 0.20 | 0.40 | 1.00 |
| c | 0.15 | 0.20 | | |
| d | 0.05 | | | |

$a \to 0$
$b \to 10$
$c \to 110$
$d \to 111$

0.5 ← 1 bit of entropy  $\frac{1}{2^1}$

0.25 ⟩ 2 bit

0.125 ⟩

0.125 ← 3 bits of entropy  $\frac{1}{2^3}$

result is the same!

Optimal only for integer entropies of symbols, ie for $\boxed{p(a_i) = \dfrac{1}{2^k}}$

---

## Problems with Huffman coding

- opimality
- requires 2 passes over input data
  (or buffer)   1st pass: determine $p(a_i)$
                2nd pass: encode

  ⟹ adaptive Huffman coding
     (assume a-priori distribution of $a_i$;
      update it)

- it is symbol based
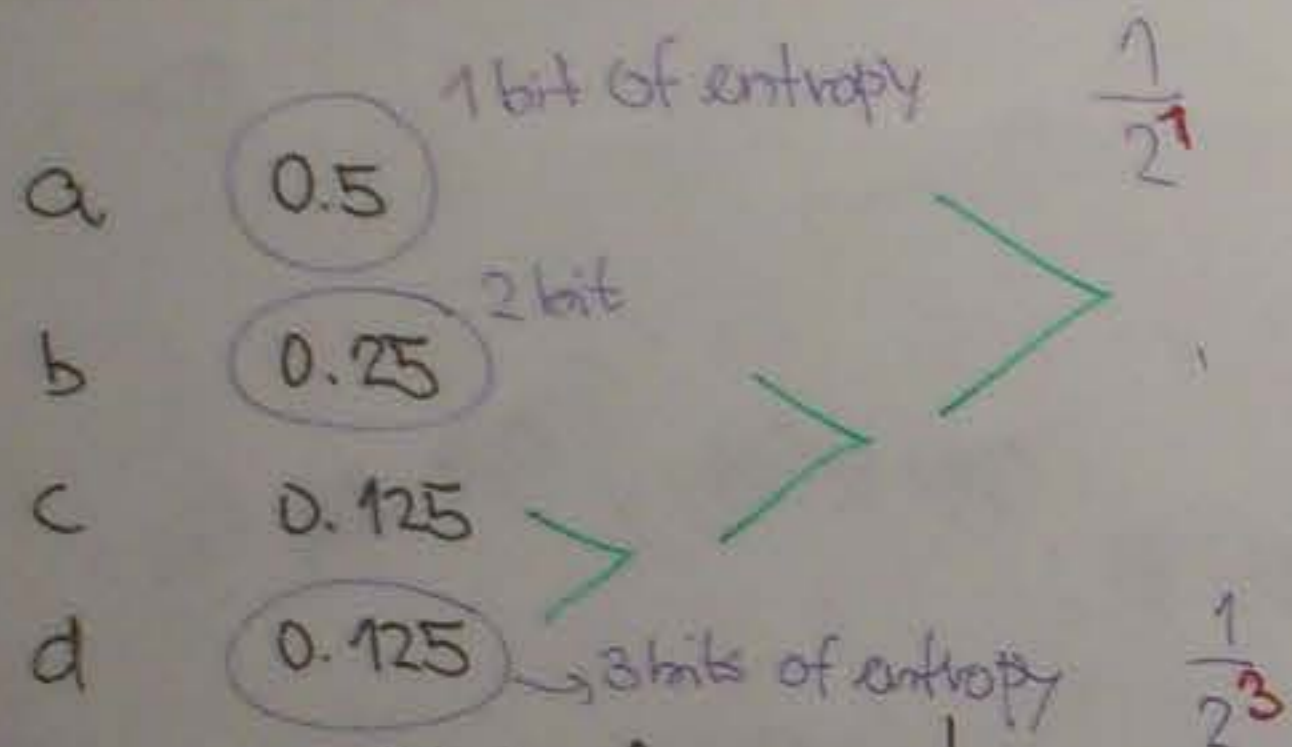  Ex: English: 'qu'

---

## Kaft's inequality

Prefix-free code with code-word lengths

$l(a_1), l(a_2), \ldots, l(a_n)$ exists iff

$$\sum_{i=1}^{\sim} 2^{-l(a_i)} \leq 1$$

Ex: $\dfrac{1}{2} + \dfrac{1}{2^2} + \dfrac{1}{2^3} + \dfrac{1}{2^3} = \dfrac{1}{2} + \dfrac{1}{4} + \dfrac{2}{8} = 1$

⟶ Full prefix-free code

Ex: $\{1, 2, 3, 3, 4\}$ ... it cannot be constructed
    $\{1, 2, 3, 4, 4\}$ ... full again

a $\quad$ (0.5) $^{1\text{ bit of entropy}}$ $\qquad \frac{1}{2^1}$

b $\quad$ (0.25) $^{2\text{ bit}}$

c $\quad$ 0.125

d $\quad$ (0.125) $\rightarrow$ 3 bits of entropy $\quad \frac{1}{2^3}$

$\rightarrow$ result is the same!

$\rightarrow$ Optimal only for integer entropies of symbols, ie for $\boxed{p(a_i) = \frac{1}{2^k}}$

## Problems with Huffman coding

- opimality
- requires 2 passes over input data
  (or buffer) $\quad$ 1st pass: determine $p(a_i)$
  $\qquad\qquad\quad$ 2nd pass: encode

$\Rightarrow$ adaptive Huffman coding
  (assume a-priori distribution of $a_i$; update it)

- it is symbol based
  Ex: English: 'qu'

## Kaft's inequality

Prefix-free code with code-word lengths $l(a_1), l(a_2), \ldots, l(a_n)$ exists iff

$$\sum_{i=1}^{N} 2^{-l(a_i)} \leq 1$$

Ex: $\frac{1}{2} + \frac{1}{2^2} + \frac{1}{2^3} + \frac{1}{2^3} = \frac{1}{2} + \frac{1}{4} + \frac{2}{8} = 1$

$\rightarrow$ full prefix-free code

Ex: $\{1,2,3,3,4\}$ ... it cannot be constructed
$\quad\;\; \{1,2,3,4,4\}$ ... full again