

CODING FOR DISCRETE SOURCES

Three major types of data sources:

- a) analog: continuous wave
- b) analog sequence: continuous in value
discrete samples
- c) discrete: symbols from finite alphabet X

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

$V = \{., -, _ \}$

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

Coding

- Source alphabet A
- Code alphabet B
- coding: symbols from $A \rightarrow$ words in B
- word: nonempty finite sequence of symbols

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

- 0 \rightarrow 00011
- 1 \rightarrow 11000
- 2 \rightarrow 10100
- 3 \vdots
- \vdots

211 \rightarrow 10100; 11000; 11000

Unique decodability

- a) $a_1 \neq a_2 \in A \rightarrow K(a_1) \neq K(a_2) \in B$
- b) no two combinations of source symbols produce the same sequence of code words

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

$K(A)$:

0	\rightarrow	0000
1	\rightarrow	0001
\vdots		
A	\rightarrow	1010
\vdots		
F	\rightarrow	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

$$M = \lceil \log_2 m \rceil \quad m \leq 2^M$$

Variable-length codes

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

$$K(a) = 0$$

$$K(b) = 10$$

$$K(c) = 11$$

- code words of different lengths
- buffering

→ Prefix-free codes

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

Ex: Morse

E → ·

T → -

A → · -

Y → - · - -

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

a → 0

b → 1

c → 10

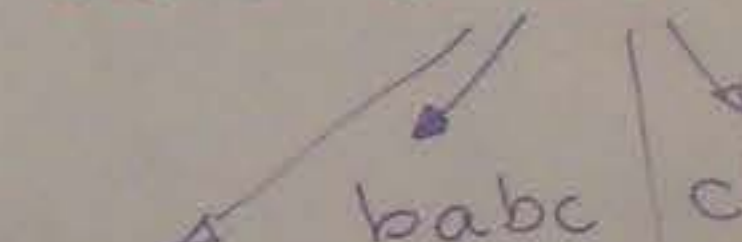
Ex:

a → 1

b → 10

c → 100

babc → 10110



babba cbba

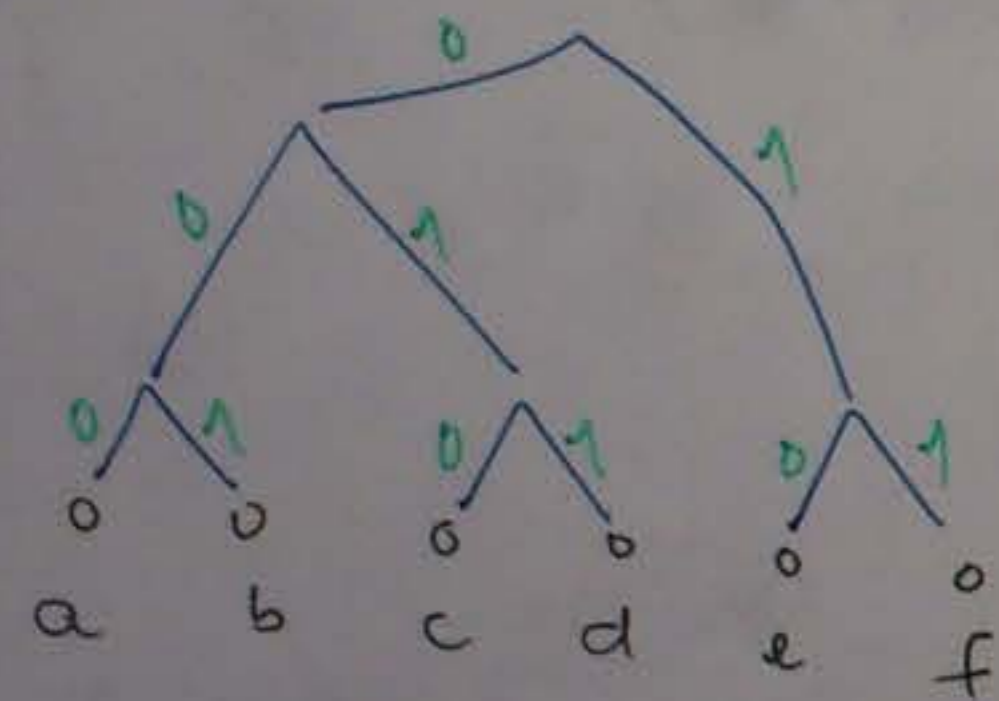
a → 1

b → 01

c → 001

Construction of a prefix-free code

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



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

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

→ all source symbols form a leaf of a binary tree

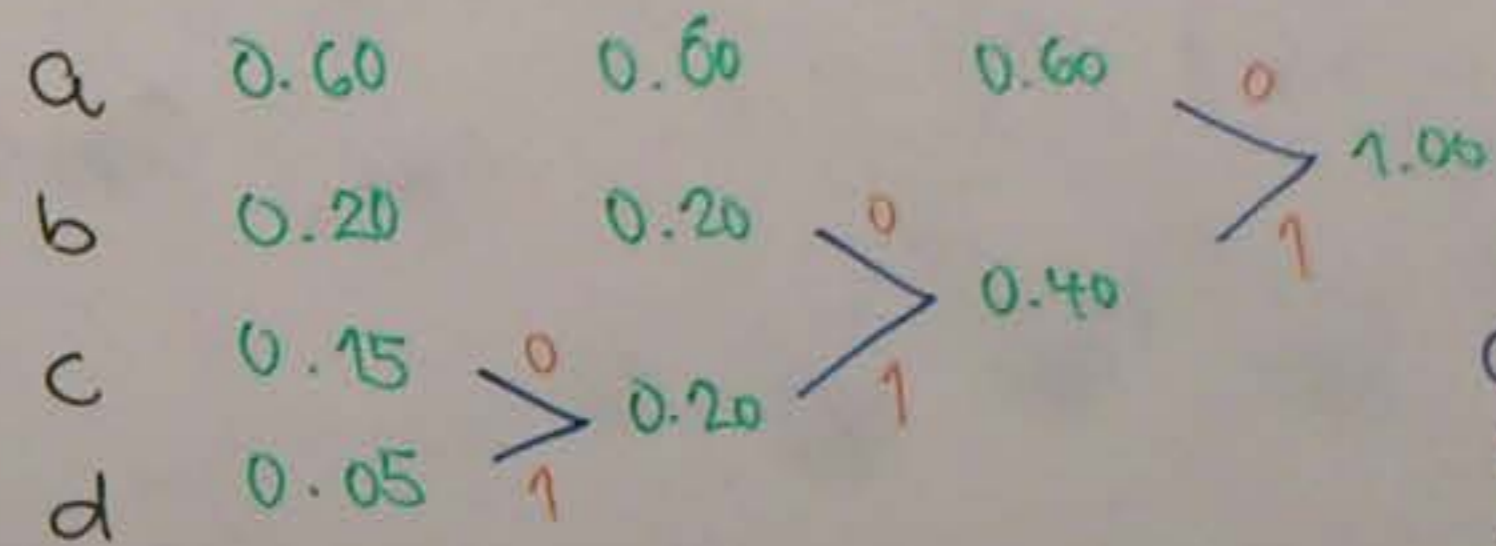
→ travelling the tree from root produces the code word

Optimal code word length

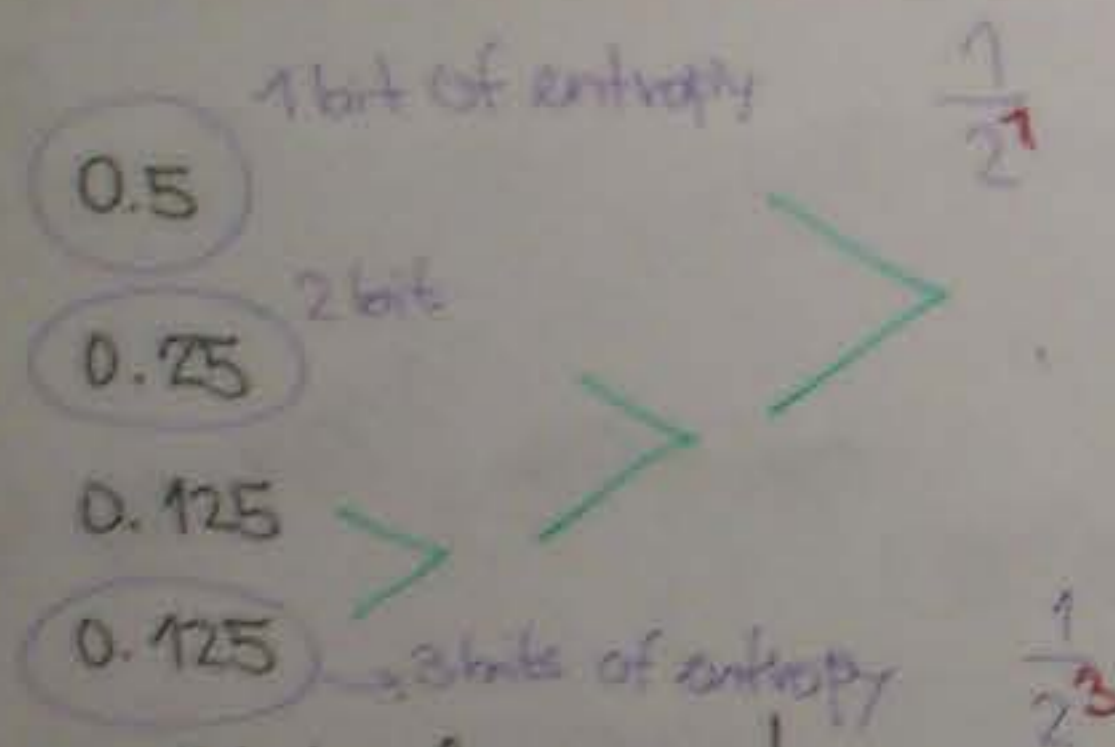
$$\mathcal{X} = \{a_1, a_2, \dots, a_m\} + \text{probabilities } p(a_1), p(a_2), \dots, p(a_m)$$

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

Huffman's algorithm



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



result is the same!
 Optimal only for integer entropies of symbols, ie for $p(a_i) = \frac{1}{2^k}$

Problems with Huffman coding

- optimality
- 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'

Kraft's inequality

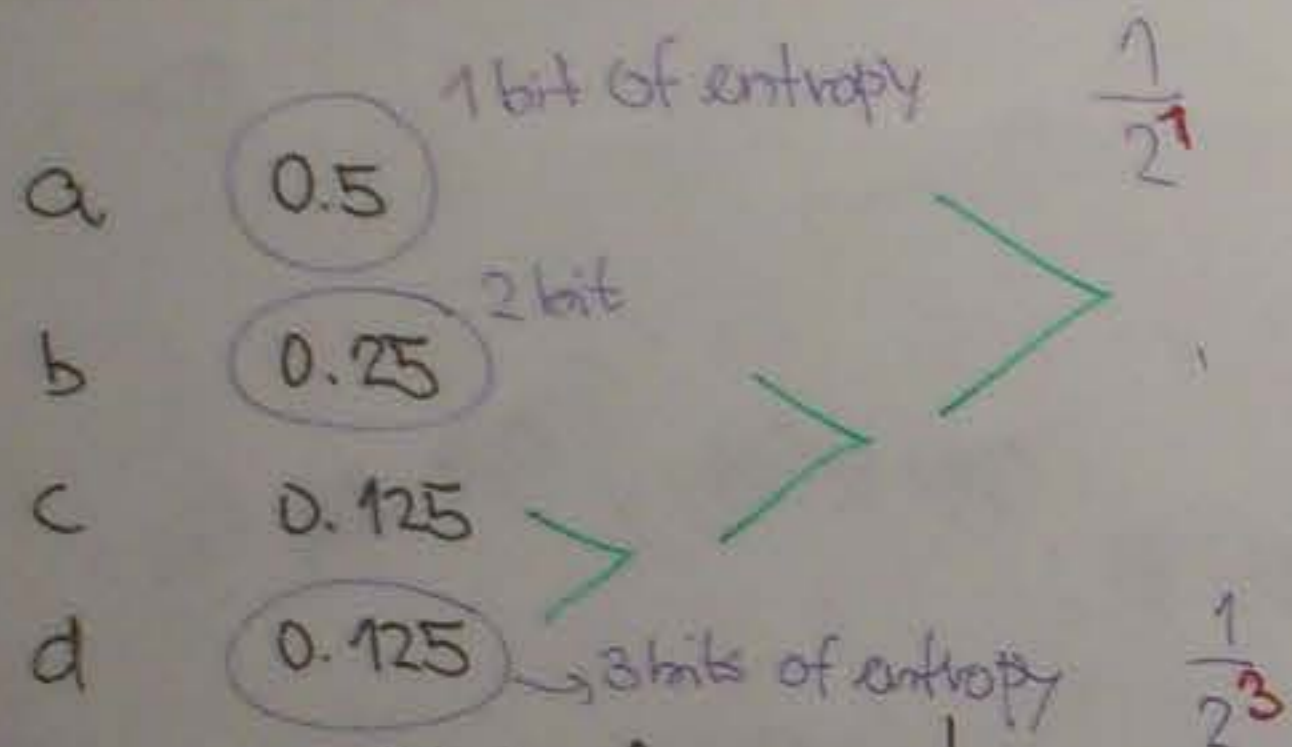
Prefix-free code with code-word lengths $l(a_1), l(a_2), \dots, l(a_n)$ exists iff

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

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

→ Full prefix-free code

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



→ result is the same!
 → Optimal only for integer entropies of symbols, ie for $p(a_i) = \frac{1}{2^k}$

Problems with Huffman coding

- optimality
- 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'

Kraft's inequality

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

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

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

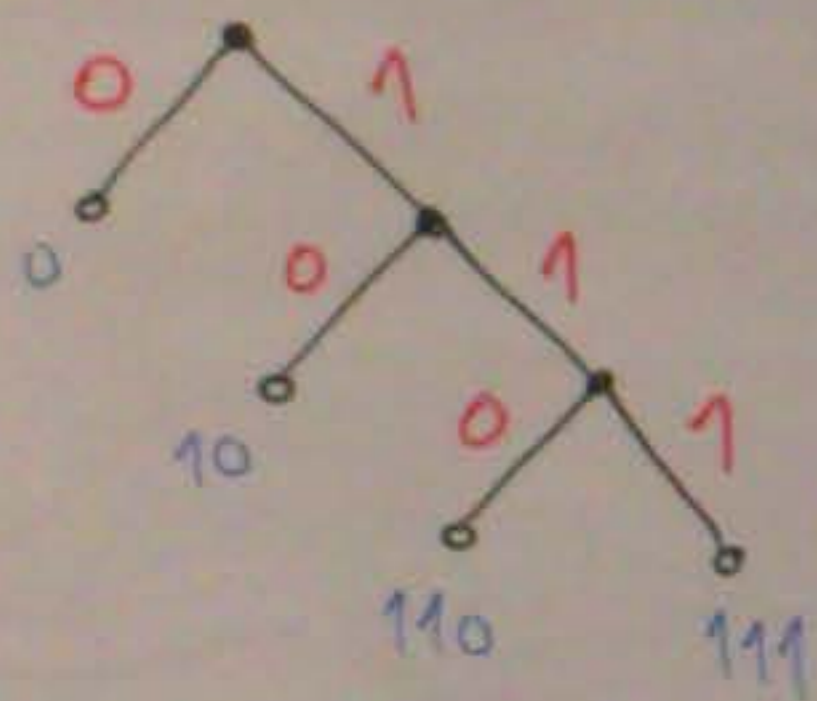
→ Full prefix-free code

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

Huffman coding

- prefix-free code
symbol length is given a-priori
(see Kraft inequality)

Ex: $l_i \in \{1, 2, 3, 3\}$



Q: why $\{1, 2, 3, 3\}$ and not $\{2, 2, 2, 2\}$?

Q: which source symbol shall be assigned when?

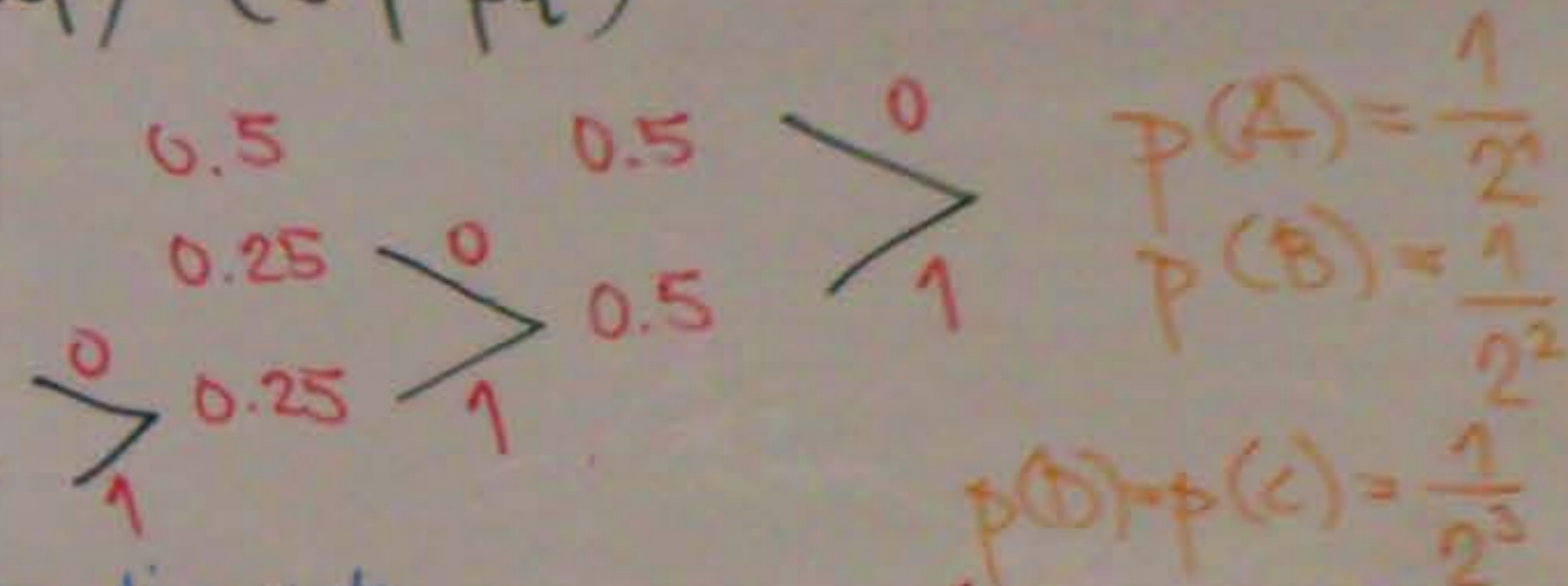
↳ this should reflect the entropy (σ, p_i)

Ex: $A \rightarrow p(A) = 0.5$

$B \rightarrow p(B) = 0.25$

$C \rightarrow p(C) = 0.125$

$D \rightarrow p(D) = 0.125$

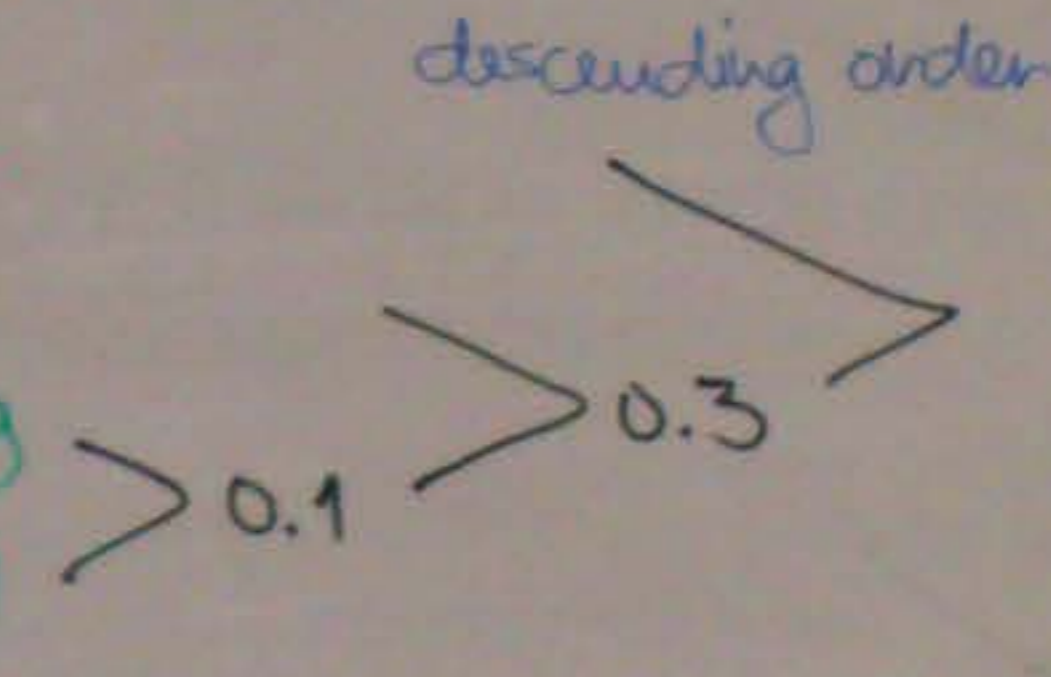


Ex: $p(A) \rightarrow 0.7$

$p(B) \rightarrow 0.2$

$p(C) \rightarrow 0.09$

$p(D) \rightarrow 0.01$



$p(D) + p(C) = \frac{1}{2^3}$

↳ optimal!

Huff. c. works well
any for these!!

Huffman coding

• problems:

a) optimality

b) 2-pass approach

pass #1: determine p_i (source symbol probabilities)
pass #2: encode

⇒ adaptive Huffman coding

start with a-priori distribution

and update it on-line

→ re-construction of Huffman tree

c) it is symbol-based

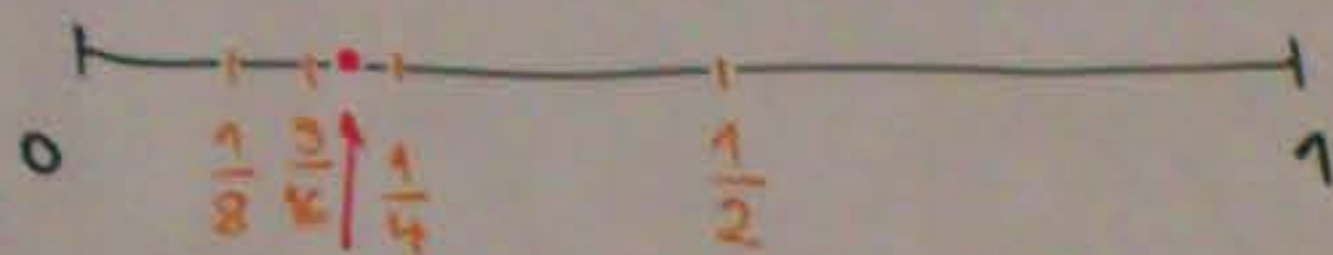
in English: 'the', 'a', 'is', 'qu'
are very frequent

Arithmetic coding

Arithmetic coding will be a topic of lab #5...

arbitrary precision arithmetic

Idea: Represent the message as a rational number and represent this rational number in binary arithmetic on $[0,1]$



$$\frac{11}{16} = \frac{1}{2} + \frac{0}{4} + \frac{1}{8} + \frac{1}{16}$$

0.1011

a) $\frac{a}{b}$... a fraction representing the message

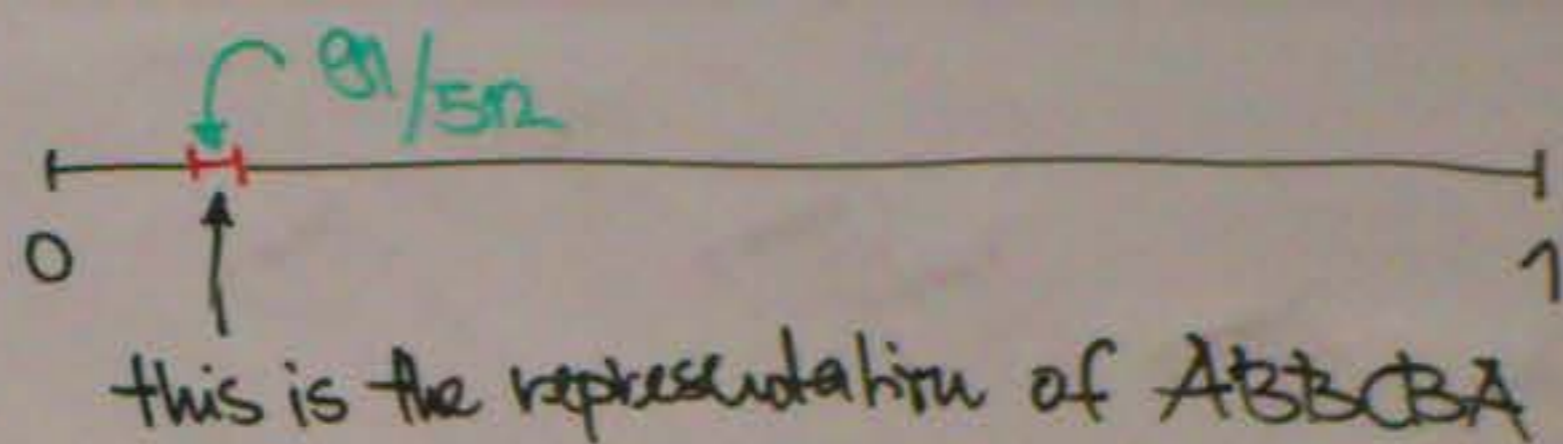
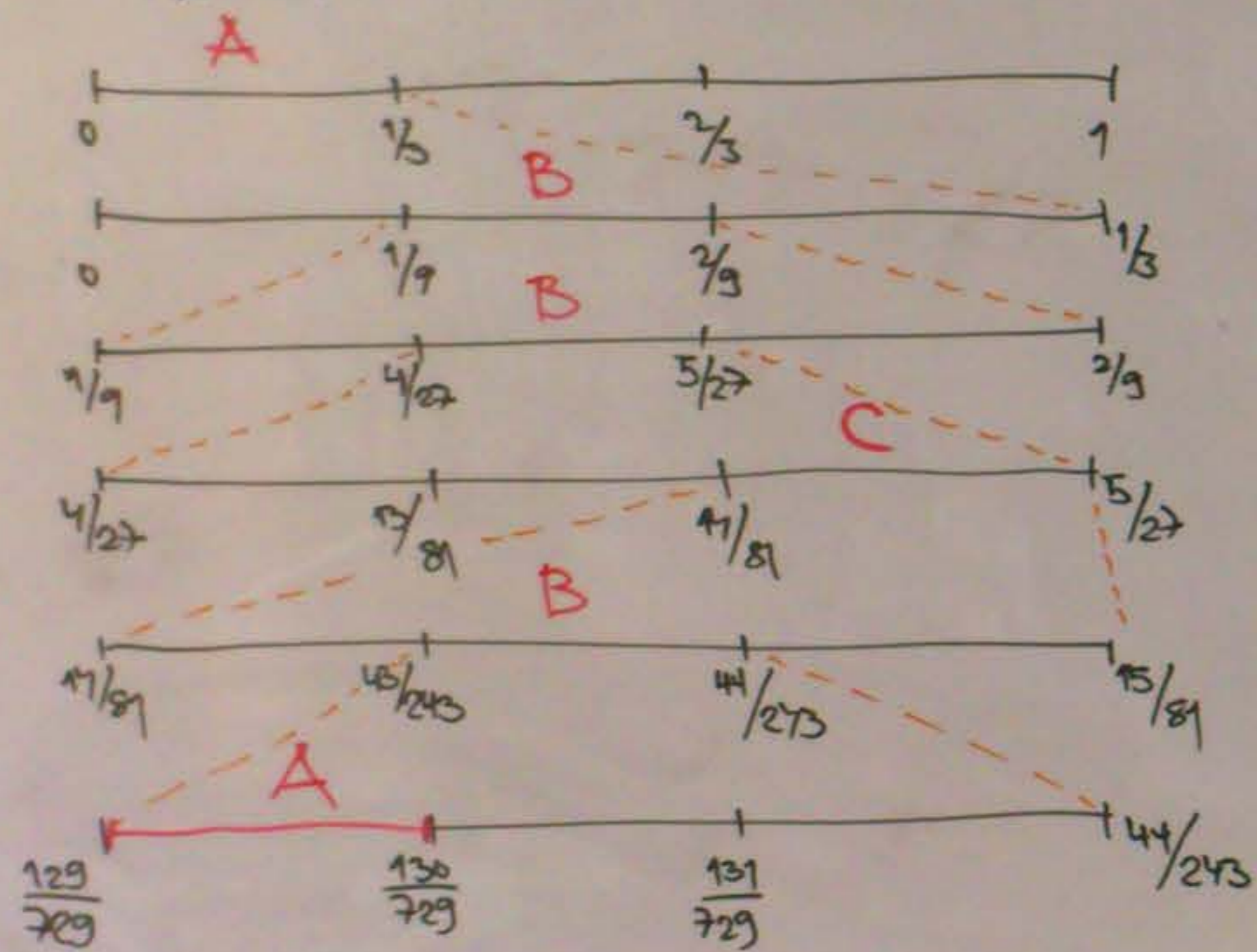
$$0.1_2 = 2^{-1} = \frac{1}{2}, \quad 0.01_2 = 2^{-2} = \frac{1}{4}$$

b) look for an interval in $[0,1]$ where $\frac{a}{b}$ lies by binary search (bisection)

→ 'bracket' Ex: $(0.1)_{10} \equiv (0.0001)_2 \in [\frac{1}{16}, \frac{2}{16}]$
and $(0.2)_{10} \in [\frac{2}{16}, \frac{3}{16}]$

Ex: $p(A) = p(B) = p(C) = 1/3$

Arithmetic code for **ABBCCA**



now: find a binary fraction $\frac{x}{2^k}$ such that

$$\frac{x}{2^k} \in \left(\frac{129}{729}, \frac{130}{729} \right)$$

! the 'shortest' one!



the binary string representing the fraction is the encoded message!

$$\frac{129}{729} = 0.17695 < \frac{x}{2^k} < 0.17832 = \frac{130}{729}$$

$$k=1: \frac{1}{2} > 0.17 \quad \times$$

$$k=2: \frac{1}{4} > 0.17 \quad \times$$

$$k=7: \frac{128}{128} < 0.176, \frac{129}{128} > 0.178 \quad \times$$

$$k=8: \frac{128}{256} < 0.176, \frac{129}{256} > 0.178 \quad \times$$

$$k=9: \frac{91}{512} > 0.176, \frac{91}{512} < 0.178 \quad \checkmark$$

binary rep. of $\frac{91}{512} = \frac{1}{8} + \frac{1}{32} + \frac{1}{64} + \frac{1}{256} + \frac{1}{512} \Rightarrow 0.001011011$