

20SK – Signals and Codes

Lecture 08 – Coding for discrete data sources, prefix-free codes (2018/11/19)

Topics discussed:

- Formal definition of a code
- Fixed-length codes for discrete sources
- Variable-length codes for discrete sources
- Unique decodability
- The Kraft inequality for prefix-free codes
- Optimal prefix-free codes, Huffman coding

The relevant literature is [1, chapter 2], [2, chapter 2] and [3, section 1.2]. Huffman coding and arithmetic coding are also described in Wikipedia.

Resources

- [1] Gallager, R.: Course materials for 6.450 *Principles of Digital Communications I*, Fall 2006. MIT OpenCourseWare (<http://ocw.mit.edu/>), Massachusetts Institute of Technology.
- [2] Adámek, J: *Foundations of Coding: Theory and Applications of Error-Correcting Codes with an Introduction to Cryptography and Information Theory*. Wiley Interscience, 1991, 352 pp.
- [3] Seibt, P.: *Algorithmic Information Theory – Mathematics of Digital Information Processing*. Springer, 2006, 447 pp.

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

$X = \{0, 1, \dots, 15\}$

$X = \{0x00, \dots, 0xffff\}$

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, \dots , 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

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

211 \rightarrow 101001100011000

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: $U = \{a, b, c\}$

$$K(a) = 0$$

$$K(b) = 10$$

$$K(c) = 11$$

- 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: $U = \{a, b, c\}$

$$a \rightarrow 0$$

$$b \rightarrow 1$$

$$c \rightarrow 10$$

babc → 10110

babca → 101101

babba → 1011011

Ex:

$$a \rightarrow 1$$

$$b \rightarrow 01$$

$$c \rightarrow 001$$

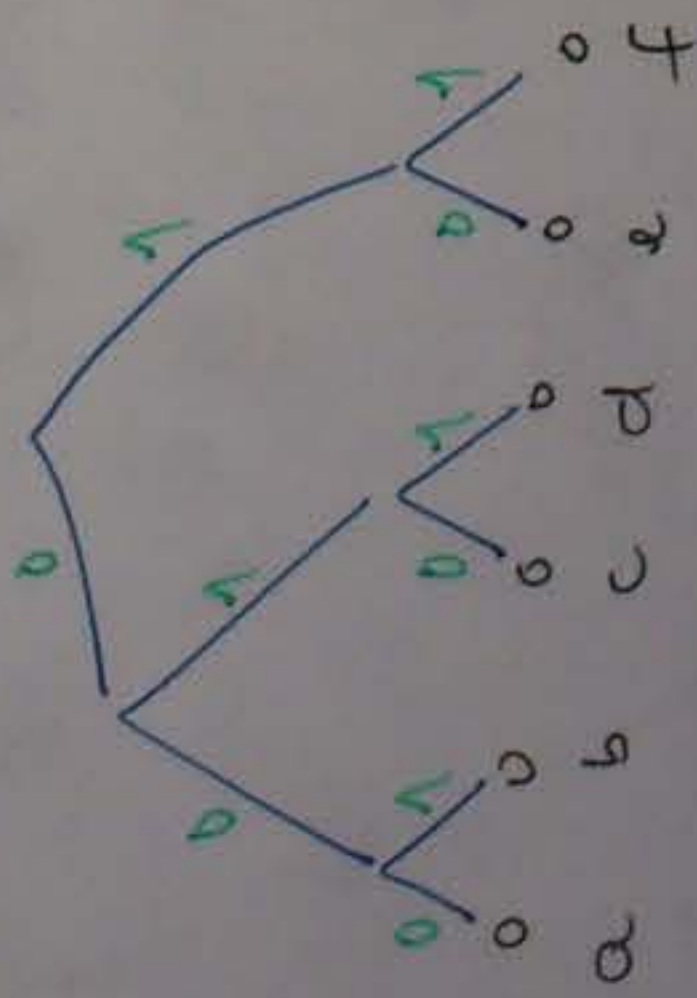
$$a \rightarrow 1$$

$$b \rightarrow 10$$

$$c \rightarrow 100$$

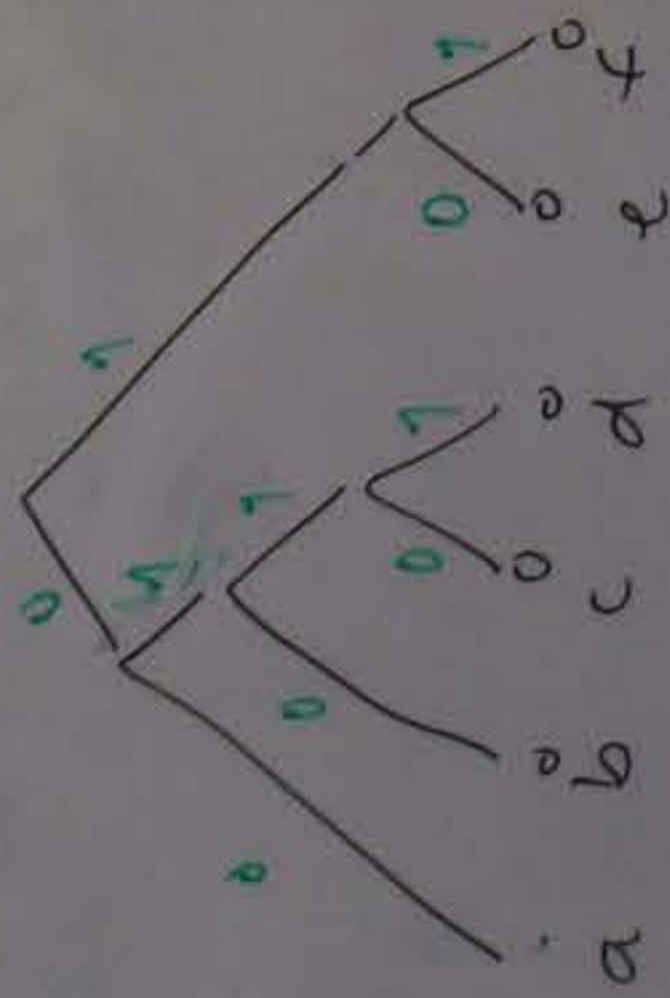
Construction of a prefix-free code

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



a → 000
 b → 001
 c → 010
 d → 011
 e → 10
 f → 11

a → 00
 b → 010
 c → 0110
 d → 0111
 e → 10
 f → 11



→ all source symbols

form a leaf of a binary tree

→ travelling the tree from root produces the code word

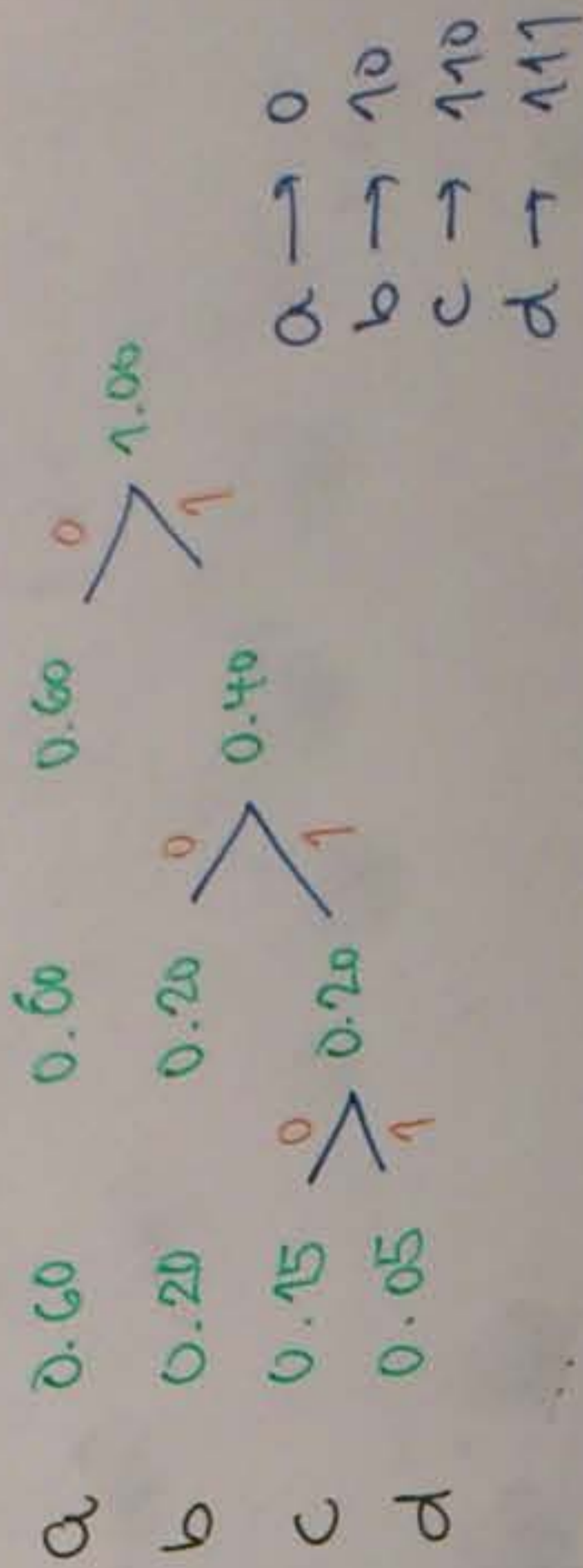
Optimal code word length

$$L = \sum_{i=1}^m a_i p_i + \text{probabilities}$$

$$P(a_1), P(a_2), \dots, P(a_m)$$

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

Huffman's algorithm



a → 0
 b → 10
 c → 110
 d → 111

1 bit of entropy

$$\frac{1}{2}$$

0.5

2 bits

0.25

0.125

0.125

3 bits of entropy

$$\frac{1}{8}$$

result is the same!

Optimal only for integer entropies of symbols, i.e. 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 $l(a_1), l(a_2), \dots, l(a_m)$ exists iff

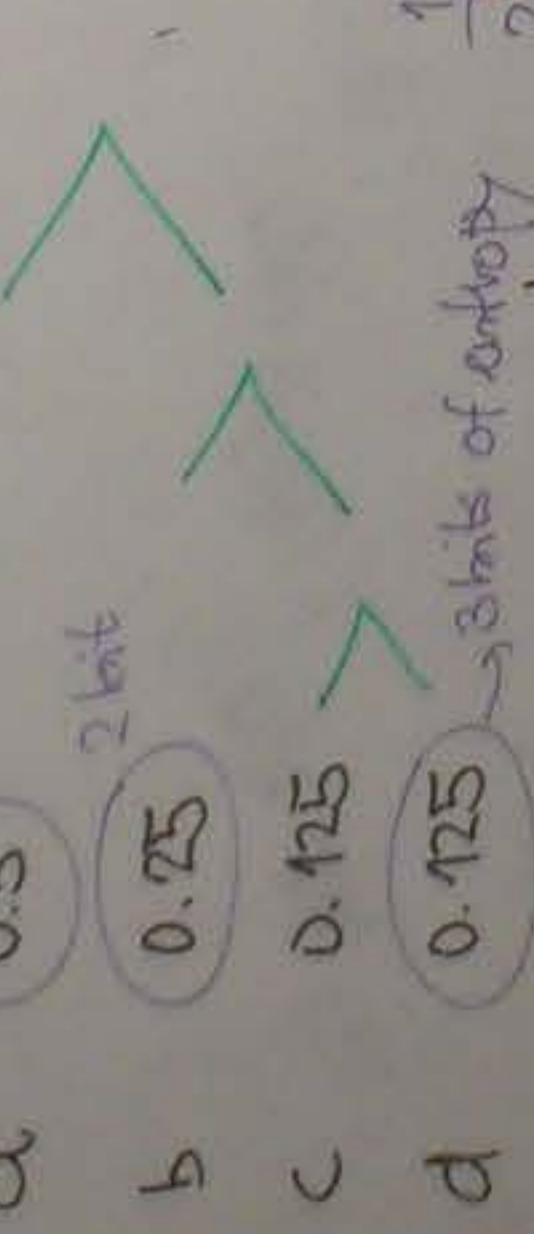
$$\sum_{i=1}^m 2^{-l(a_i)} \leq 1$$

Ex: $\frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = 2 > 1$
 $\frac{1}{2} + \frac{1}{2^2} + \frac{1}{2^2} + \frac{1}{2^2} + \frac{1}{2^2} = 1$

→ Full prefix-free code

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

1 bit of entropy $\frac{1}{2}$



→ result is the same!

→ Optimal only for integer entropies of symbols, i.e. 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: $\frac{1}{2} + \frac{1}{2^2} + \frac{1}{2^3} + \frac{1}{2^3} = \frac{1}{2} + \frac{1}{4} + \frac{2}{8} = 1$

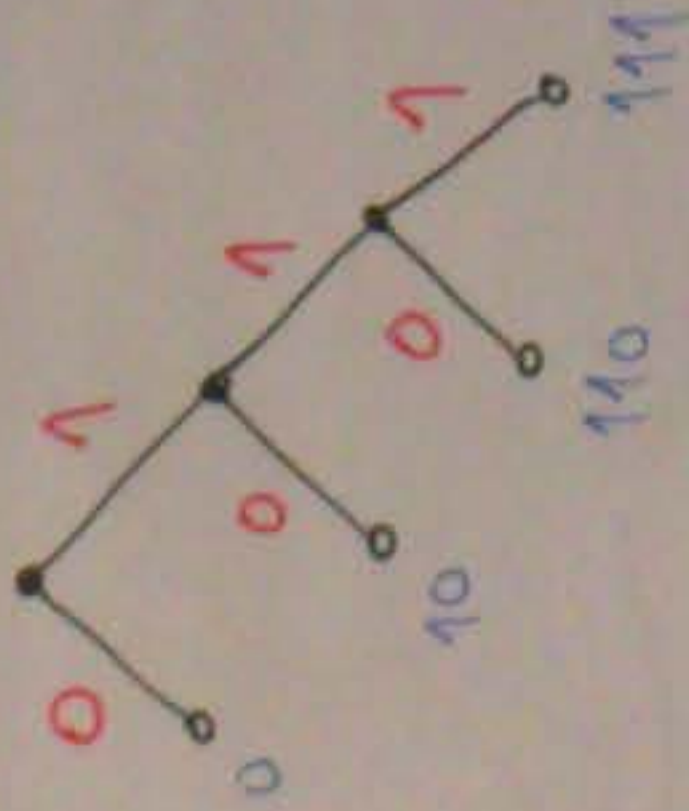
→ Full prefix-free code

Ex: $\{1, 2, 3, 4\}$... 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: $\mathcal{L}_i \in \{1, 2, 3, 1, 3\}$

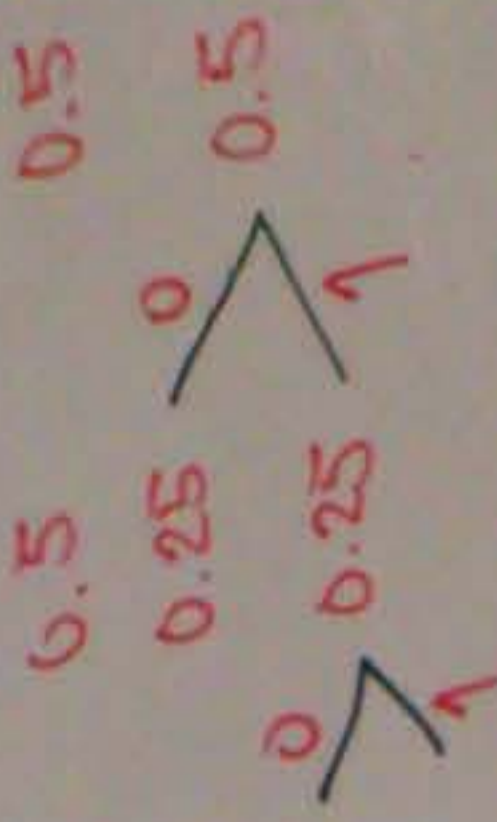


Q: why $\{1, 2, 3, 1, 3\}$ and not $\{2, 2, 2, 2, 2\}$?

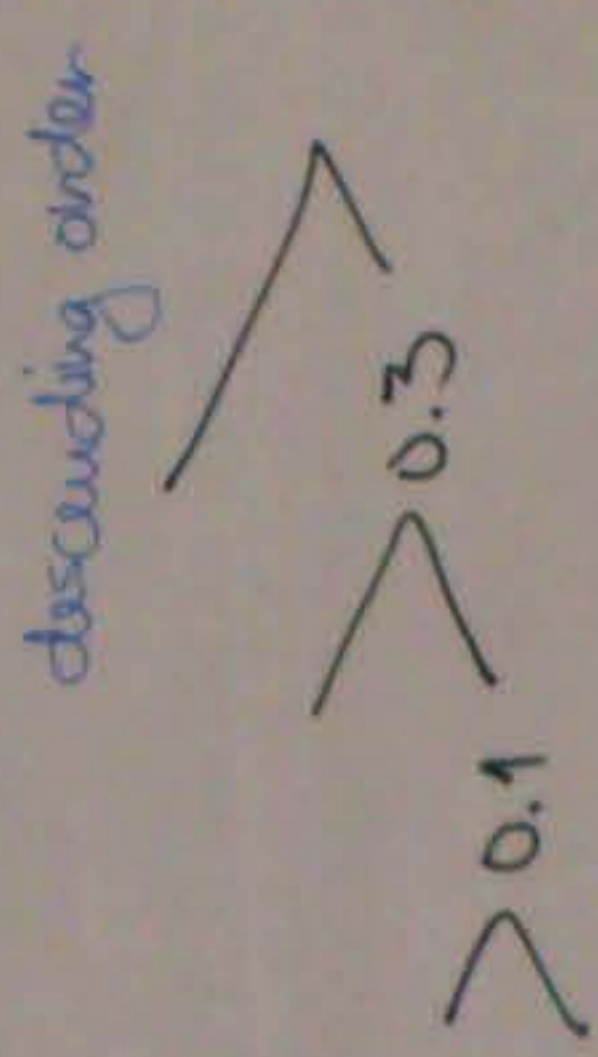
Q: which source symbol shall be assigned when?

→ this should reflect the entropy (α_i, p_i)

Ex: $A \rightarrow p(A) = 0.5$
 $B \rightarrow p(A) = 0.25$
 $C \rightarrow p(C) = 0.125$
 $D \rightarrow p(D) = 0.125$



Ex: $p(A) \rightarrow 0.7$
 $p(B) \rightarrow 0.2$
 $p(C) \rightarrow 0.09$
 $p(D) \rightarrow 0.01$



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

→ optimal!
 Huff. c. works well
 only for these!!

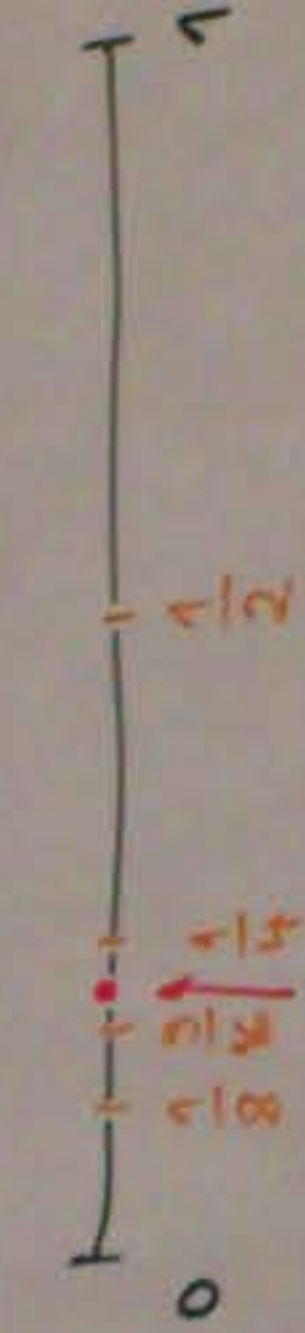
Huffman coding

- problems:
 - optimality
 - 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', 'gu'
are very frequent

Arithmetic coding

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

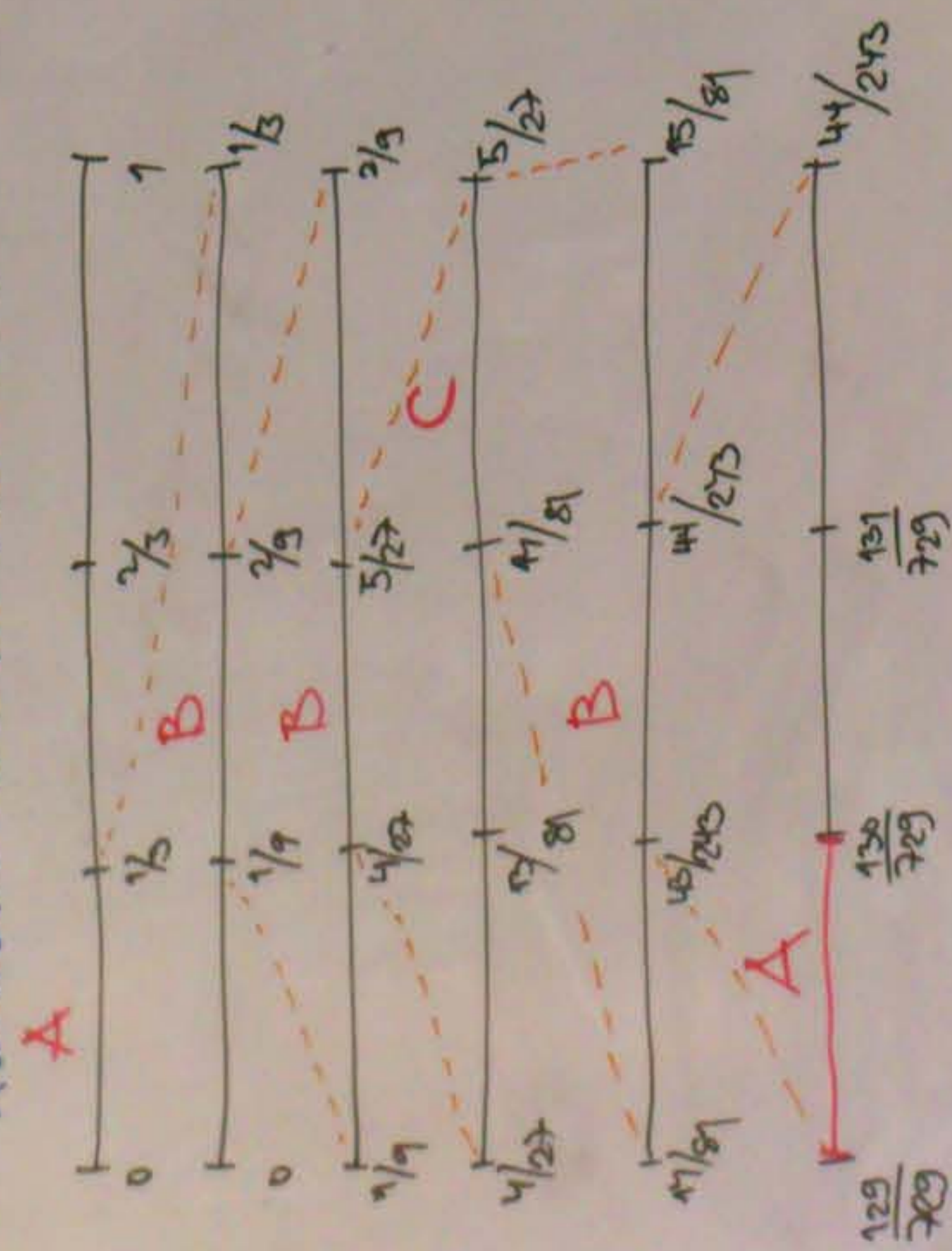
$$0.1_2 = \bar{2} = \frac{1}{2}, \quad 0.01_2 = \bar{2}^2 = \frac{1}{4}$$

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 ABCBA



this is the representation of ABCBA

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 binary string representing the fraction is the encoded message!

$$\frac{129}{729} = 0.17635 < \frac{x}{2^k} < 0.17832 = \frac{130}{729}$$

k=1: $\frac{1}{2} > 0.17$ X
 k=2: $\frac{1}{4} > 0.17$ X
 ...

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

! the 'shortest' one!

k=7: $\frac{22}{128} < 0.176$ | $\frac{23}{128} > 0.178$ X
 k=8: $\frac{45}{256} < 0.176$ | $\frac{46}{256} > 0.178$ X
 k=9: $\frac{91}{512} > 0.176$ | $\frac{91}{512} < 0.178$ ✓

$\Rightarrow 0.001011011$