

# 20SK – Signals and Codes

---

## Lecture 09 – Dictionary coders (2018/11/26)

Topics discussed:

- Dictionary-based compression schemes in general
- RLE
- Lempel-Ziv (LZ77) source coding
- LZW

The relevant literature is [1,2] for LZ77 and [3] for LZW. For more insight into data compression methods see also books by David Salomon [4,5]. Some relevant information can be also found in Wikipedia and other online sources.

### Resources

- [1] Shor, P.: *Lempel-Zip notes* [online]. MIT, 2005. Available from: [http://www-math.mit.edu/~shor/PAM/lempel\\_ziv\\_notes.pdf](http://www-math.mit.edu/~shor/PAM/lempel_ziv_notes.pdf)
- [2] Ziv, J., Lempel, A.: *A Universal Algorithm for Sequential Data Compression*. IEEE Transactions on Information Theory 23 (3), 1977, pp. 337–343. doi:10.1109/TIT.1977.1055714.
- [3] Welch, T.: *A Technique for High-Performance Data Compression*. Computer 17 (6), 1984, pp. 8–19. doi: 10.1109/MC.1984.1659158.
- [4] Salomon, D.: *Concise Introduction to Data Compression*. London: Springer Verlag, 2008, 311 pp.
- [5] Salomon D., Motta G.: *Handbook of Data Compression*. London: Springer Verlag, 2010, 1360 pp. doi: 10.1007/10.1007/978-1-84882-903-9

## Huffman Coding

A ...  $p_A = \frac{1}{3}$   
 B ...  $p_B = \frac{1}{3}$   
 C ...  $p_C = \frac{1}{3}$

$$H(X) = -\log_2 \frac{1}{3} = 1.585$$

Example: Huffman code

A  $\rightarrow$  0  
 B  $\rightarrow$  10  
 C  $\rightarrow$  11

ABBCBA  $\rightarrow$  0101011100  
 $\rightarrow$  10 bits, 6 symbols

$$E(H(X)) = \frac{10}{6} = 1.67 \text{ bits/symbol}$$

ABBCB  $\rightarrow$  01010110

$\Rightarrow$  9 bits, 5 symbols

$$E(H(X)) = \frac{9}{5} = 1.8 \text{ bits/symbol}$$

ABCABC  $\rightarrow$  0101101011

$$E(H(X)) = \frac{10}{6} = 1.67 \text{ bits/symbol}$$

works well only for  $p_i = 2^{-k}$   
 bad for  $p_i$ 's like 75%, 37.5%, ...  
 $\frac{3}{4}, \frac{5}{8}, \frac{7}{8}$

## Arithmetic Coding

idea: represent the source as real number from  $(0,1)$   
 by successive subdivision of interval  $(0,1)$  according to symbol probabilities

Example: Arithmetic code for ABBCBA



$\Rightarrow$  Combination of binary fractions  
 $\left(\frac{a_i}{2^i}\right)$  so that  $\sum_{i=0}^{\infty} \frac{a_i}{2^i} \in \left(\frac{129}{729}, \frac{130}{729}\right)$



continued....

$$\frac{129}{729} < m < \frac{130}{729}$$

$$\uparrow \frac{a_i}{2^i}$$

$$0.17695 < m < 0.17852$$

$i=7$  ... no match for  $m$  :

$i=8$  ... no match for  $m$  :

$$i=9 \dots m = \frac{91}{512} - \frac{1}{512} + \frac{1}{256} + \frac{1}{64} + \frac{1}{32} + \frac{1}{8}$$

$$m = 0.1001011011$$

ABBCBA  $\rightarrow$  00101011

$\Rightarrow$  9 bits, 6 chars

$$E(M(x)) = \frac{9}{6} = 1.5 \text{ bits/symbol}$$

1) arbitrary precision arithmetic

2) needs to read the whole message beforehand

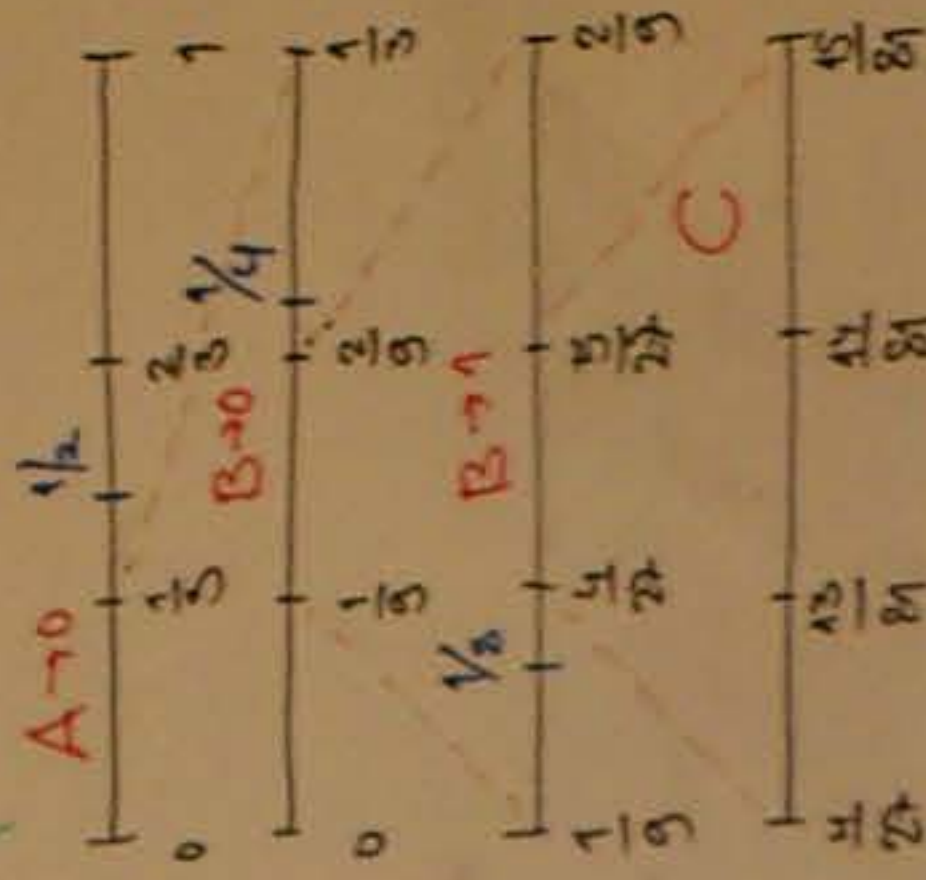
$\Rightarrow$  better algorithms exist for efficient computation of  $m$

$$\frac{22}{128} < 0.176 \text{ and } \frac{23}{128} > 0.179$$

$$\frac{45}{256} < 0.176$$

idea: represent the source as real number from  $(0,1)$  by successive subdivision of interval  $(0,1)$  according to symbol probabilities

Example: Arithmetic code for ABBCBA



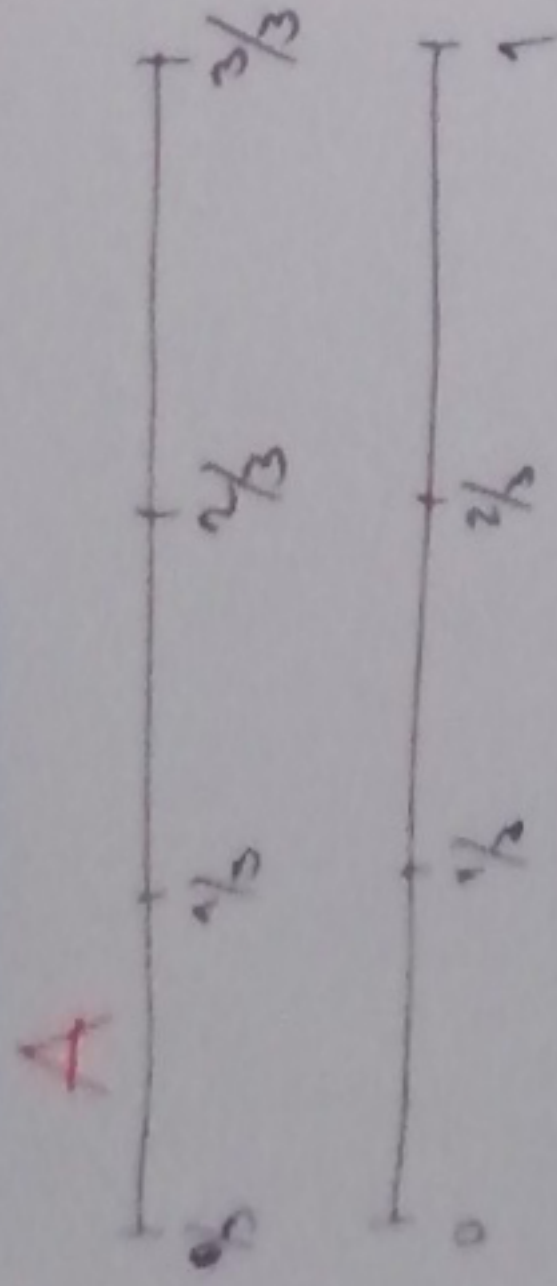
$\Rightarrow$  Combination of binary fractions  $\sum_{i=1}^n \frac{a_i}{2^i} < \frac{130}{729}$  so that

Example: Arithmetic code for ABBCBA



# ARITHMETIC CODING

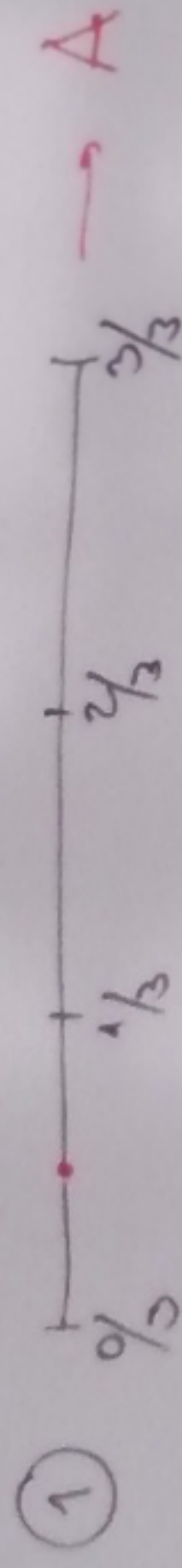
- re-normalisation



We will always use  $\langle 0,1 \rangle$  to assign intervals for letters and scale afterwards.

# DECODING

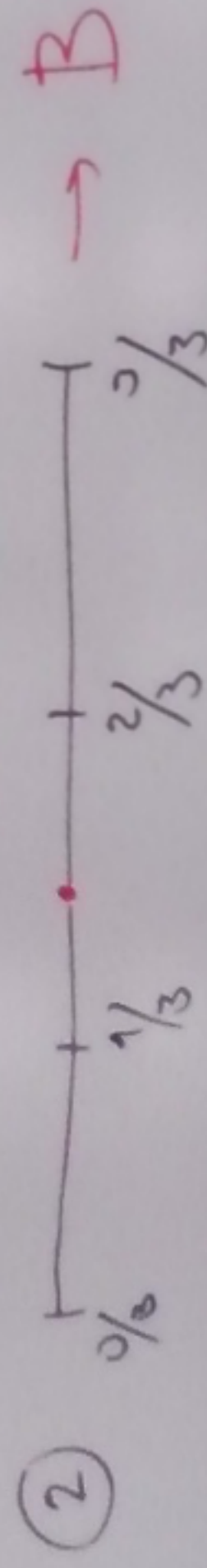
message: 0.001011011 =  $\frac{91}{512}$  0.177...



(  $\frac{91}{512}$  - lower bound of interval ) / length of interval

$$\Rightarrow \text{a number from } \langle 0,1 \rangle = \frac{273}{512}$$

$$\frac{91}{512} / \frac{1}{3} = \frac{3 \cdot 91}{512} = \frac{273}{512}$$



$$\left( \frac{273}{512} - \frac{1}{3} \right) \cdot 3 = \frac{3 \cdot 273 - 512}{512} = \frac{307}{512}$$



## Dictionary methods

- encode phrases, not single symbols into code words
- output words/symbols are equiprobable

Ex: Encoding English text

'the' → 1  
'aw' → 2  
'is' → 3  
⋮

## based on 'transformation table' → DICTIONARY

- table is dynamic → coding improves with the number of symbols processed (usually!)

## Historical remark: RLE (Run-Length Encoding)

PCX (i.i.r.c. max 16 colors)

Ex: binary B&W bitwap

0... black  
1... white

111100000101111  
4,1,6,10,11,1,10,14,17



8-bit packed

3 bits ⇒ 8 colours

5 bits ⇒ max. 32 repetitions

can be implemented as a dictionary, BUT it is computationally ineffect.

## Sliding window

- given by its width in 'elements' (source symbols)
- moves through the data set

Ex: sliding window with 8

0110101010101010

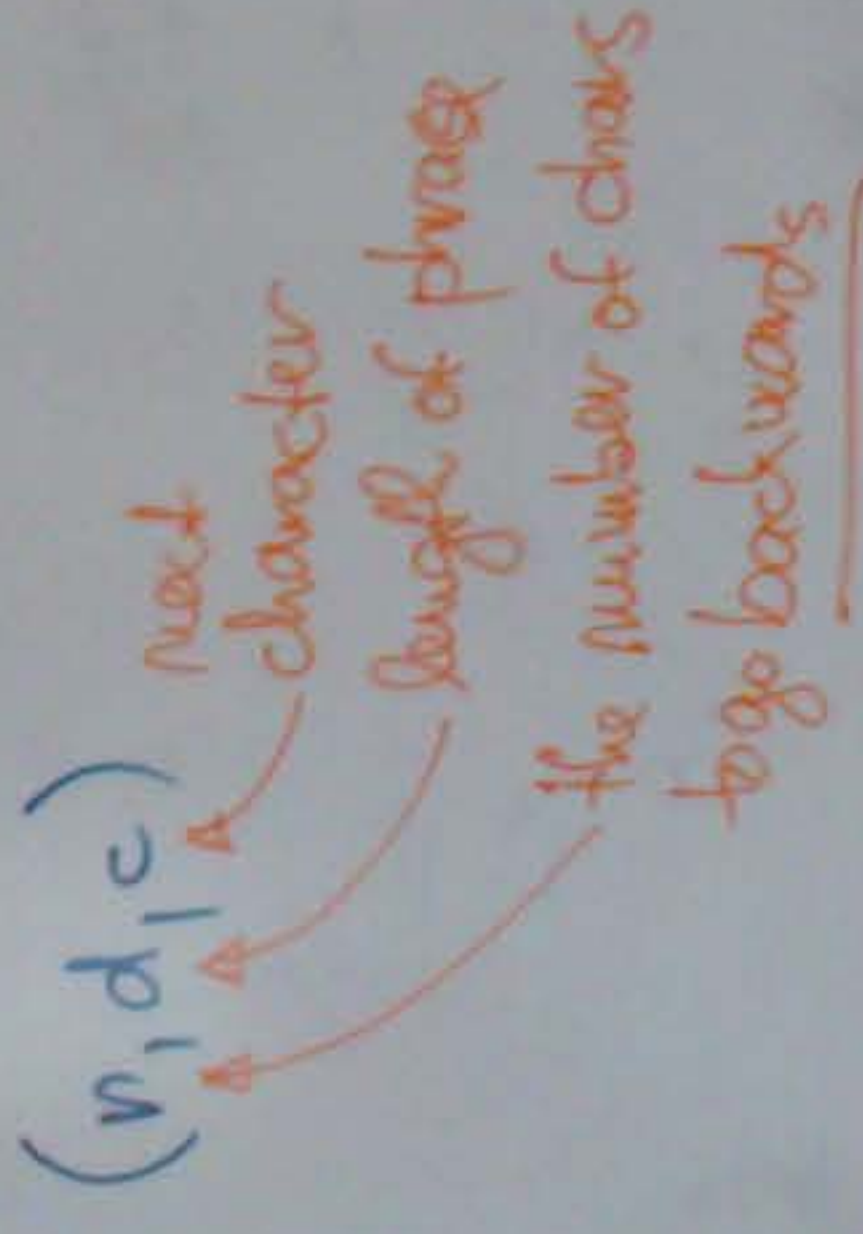
position 0

position 8



# LZ77 (Lempel-Ziv)

- use a sliding window to search for repetitions of a phrase; if found, the current phrase is replaced with a reference



K d y o y w b y l y w n y o y a w b y w c h y o y .

K d y o y w b y l y w n y o y a w b y w c h y o y .

(00,1,K) (00,1,d) (00,1,y) (00,1,w) (2,1,1,w) (3,2,1,e) (5,2,1,r) (10,4,1,n)

no previous reference found

e b y l y w b y w c h y o y .

(0,0,1,e) (12,15,1,b) (3,2,1,c) (0,0,1,h) (16,13,1,o)

32 characters → 13 triplets

x 8bit → 256 bits

x 7bit → 224 bits

(uncompressed, but

"optimized" sequence)

→ ASCII

208 bits

if 16 bits/triplet

(5bit, 4bit, 7bit) → 16 bits/triplet

max. length of phrase 15 char

length of the window is 32 char

(i.e. the max. backwards step is 32)



# LZW (Lempel-Ziv-Welch)

- search in a long sliding window is slow
- use a tree-based dictionary

a) dictionary is adaptive  
(can be discarded after compression and re-created when decompressing)

b) pre-defined basis

K d y b y - by i y r y o y n a ch y o y

K	d	y	b	y	-	by	i	y	r	y	o	y	n	a	ch	y	o	y
0	1	2	3	2	4	15	5	16	6	14	16	7	8					
+	+	+	+	+	+	....												

this goes to dictionary

-	c	h	y	o	y	·	EOF
4	9	10	22	11			

32 characters  
↳ 256 bits  
224 bits  
29 symbols à 6 bits = 132 bits

## Dictionary:

k → 0	d → 1	y → 2	b → 3	y → 4	l → 5	r → 6	n → 7	e → 8	c → 9	h → 10	.	→ 11	(basis)
kd → 12	dy → 13	yo → 14	by → 15	yr → 16	bo → 17	byl → 18	ly → 19	yur → 20	ry → 21	yby → 22	yun → 23	ne → 24	ob → 25
byy → 26	yub → 27	oy → 28	uc → 29	ch → 30	by → 31	yoy → 32							