

# Huffman Codes and Text Compression

Representing letters  $\rightarrow$  A | B | C | D |  $\emptyset$  if we represent these with binary code how many digits should we use?

ABC DAB CBBA = m

$\underbrace{000 \ 001 \ 010 \ 011 \ 000}_{A \ B \ C \ D \ A} \dots = 3 \times m$  bits

↳ encoding table

$$2^m > 5$$

use 3 digits to represent 5 chars

Variable Sized Encoding  $\rightarrow$

d = number of chars

n = number of items

$$\text{total comp} = n + d \log d$$

A  $\rightarrow$  0

B  $\rightarrow$  10  $\rightarrow$  0 ile baslayamaz

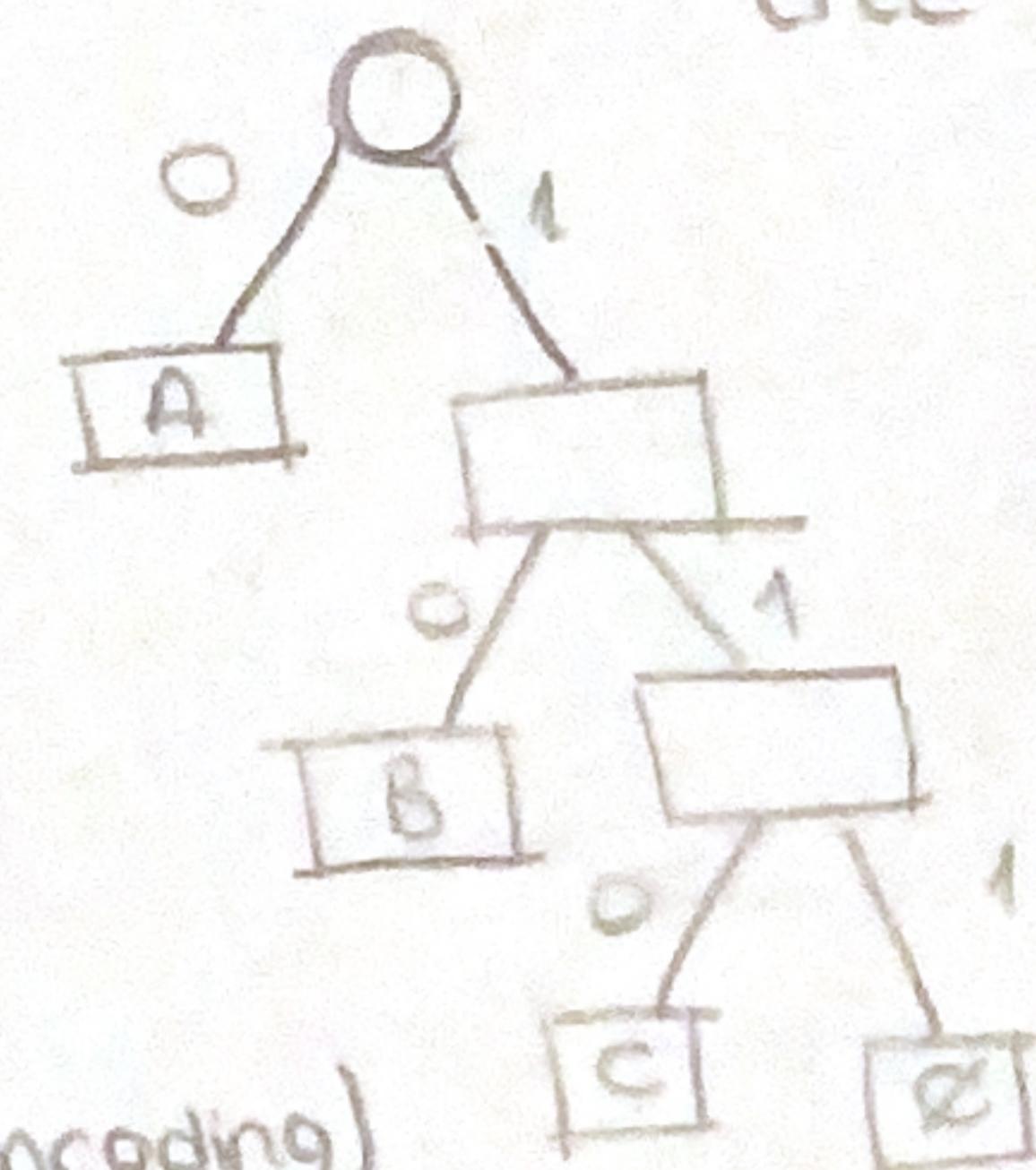
C  $\rightarrow$  110  $\rightarrow$  10 ve 0 ile baslayamaz

D  $\rightarrow$  111  $\rightarrow$  10, 110, 0 ile baslayamaz

A B C C A A B A  
0 1 0 1 1 0 1 1 0 0 1 0 0

↳ encoding tree

ex u A | B | C | D | - } letter frequencies  
0.35 | 0.1 | 0.2 | 0.2 | 0.15 }



Fixed Encoding we need 3 bits  
100 letter  $\rightarrow$   $3 \times 100 = 300$  bits  
text

If we can build up a variable encoding tree we can decrease the space (Huffman encoding)  
- En türkük frekansları recursive olarak birleştir.

<table border="1"><tr><td>A</td><td>0.35</td></tr></table>	A	0.35	<table border="1"><tr><td>B</td><td>0.1</td></tr></table>	B	0.1	<table border="1"><tr><td>C</td><td>0.2</td></tr></table>	C	0.2	<table border="1"><tr><td>D</td><td>0.2</td></tr></table>	D	0.2	<table border="1"><tr><td><math>\emptyset</math></td><td>0.15</td></tr></table>	$\emptyset$	0.15
A	0.35													
B	0.1													
C	0.2													
D	0.2													
$\emptyset$	0.15													

\* complexity of construction  $d \log d$

A  $\rightarrow$  11

B  $\rightarrow$  100

C  $\rightarrow$  00

D  $\rightarrow$  01

$\emptyset \rightarrow$  101

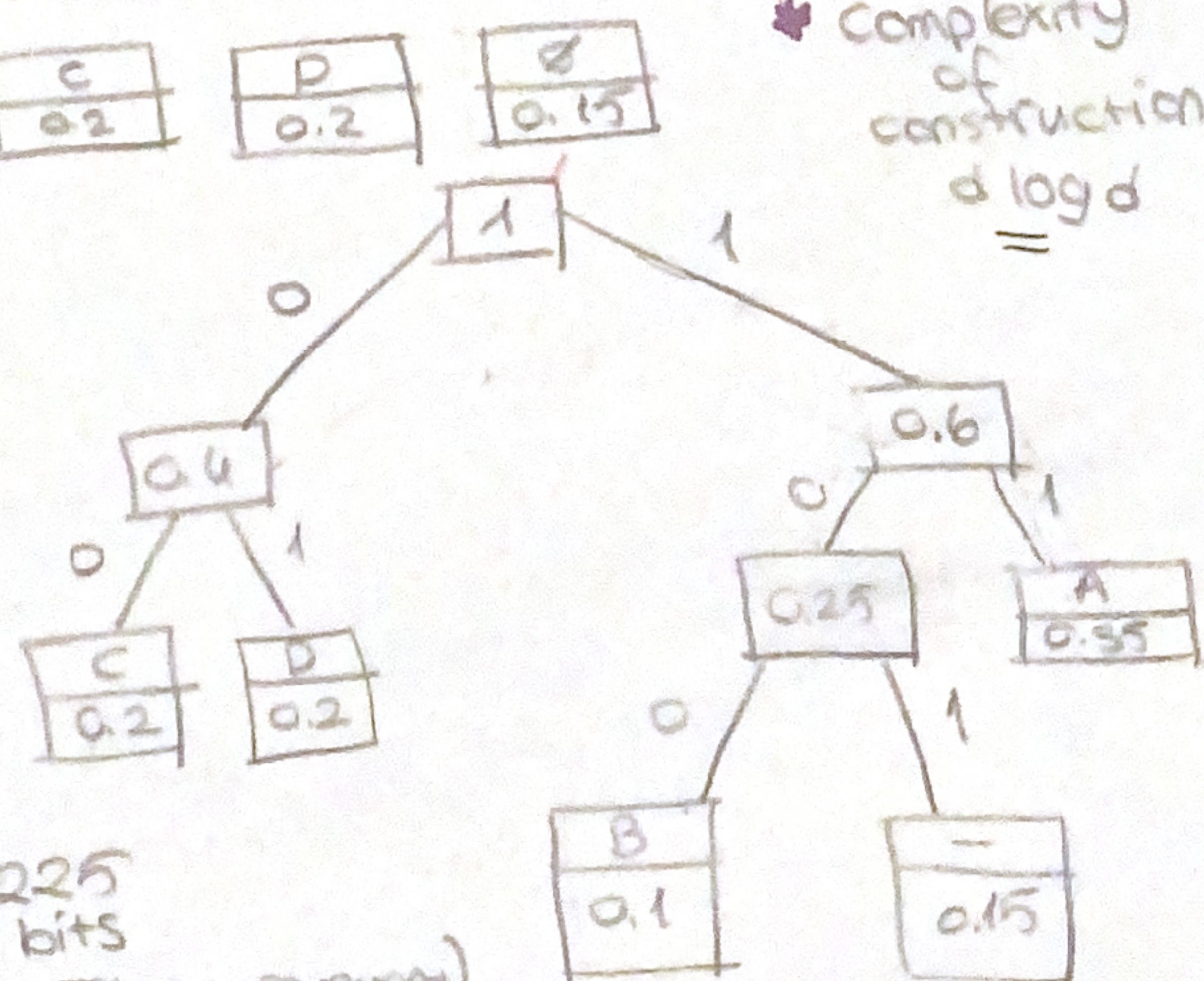
$$100 \text{ letter} \rightarrow 35 \times 2 \text{ A}$$

$$10 \times 3 \text{ B}$$

$$20 \times 2 \text{ C}$$

$$20 \times 2 \text{ D} \rightarrow 225 \text{ bits}$$

$$15 \times 3 \emptyset$$



(25% improvement)

## Huffman's Algorithm

1. Initialize  $n$  one node trees and label them with letters and sequences frequencies
2. Repeat this step until every node is combined in a single tree. Combine them and make left & right subtree of a root which has the combined frequency.

ex:

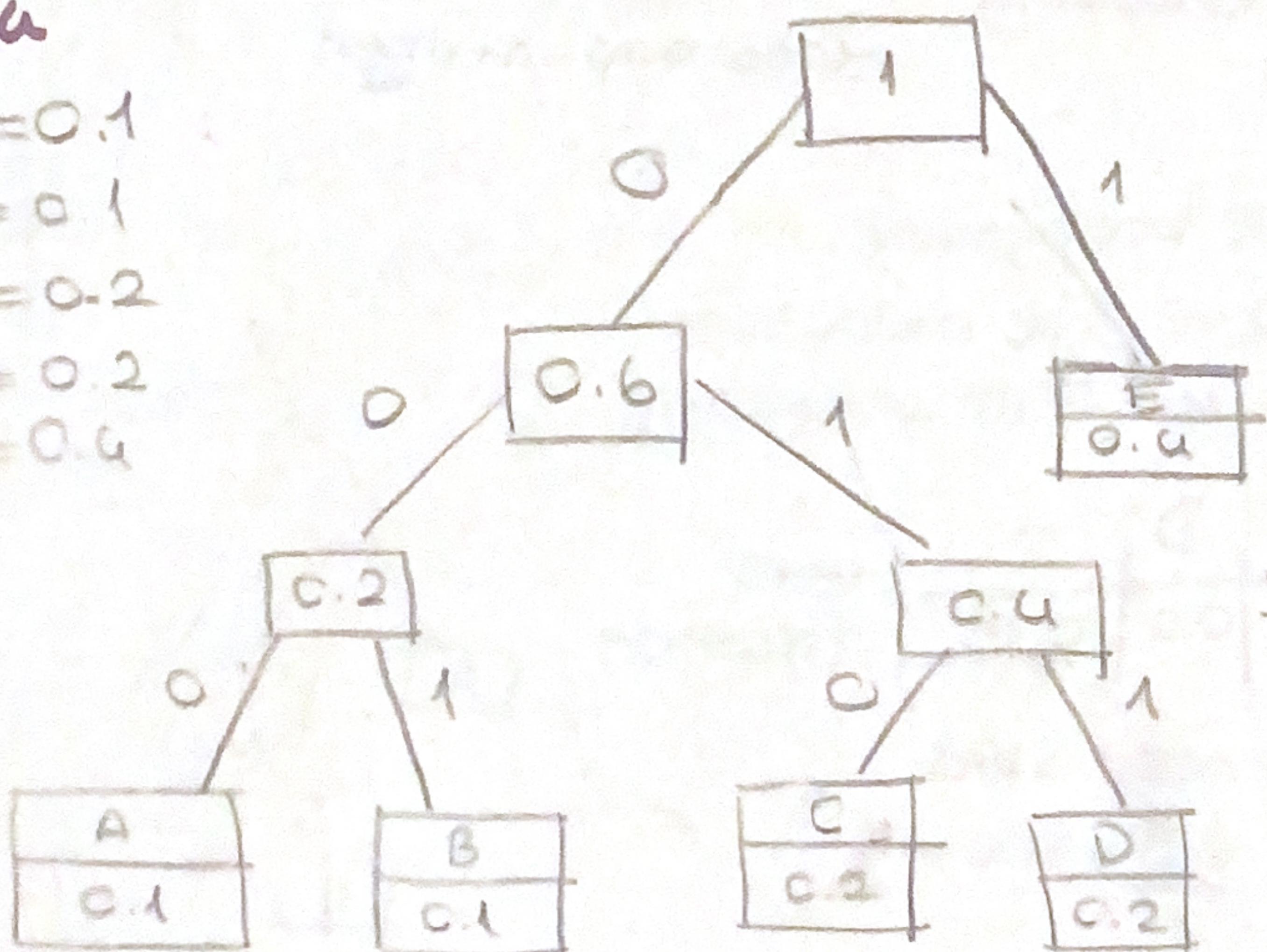
$$A = 0.1$$

$$B = 0.1$$

$$C = 0.2$$

$$D = 0.2$$

$$E = 0.4$$



$$A = 000$$

$$B = 001$$

$$C = 010$$

$$D = 011$$

$$E = 1$$

100 letters

$$A = 3 \times 10$$

$$B = 3 \times 10$$

$$C = 3 \times 20$$

$$D = 3 \times 20$$

$$E = 1 \times 40$$

Fixed Encoding  $\rightarrow 3 \times 100 = 300$  bits

Variable Encoding  $\rightarrow 220$  bits

220 bits