Problems 4 - Character codings, Dictionary coders

Coding Technology

Illés Horváth

2025/10/31

Reminder: character encodings

Character codings: prefix codes, code tree for decoding.

Fixed length encoding.

Shannon–Fano code: codeword lengths $\ell_i = \lceil -\log_2 p_i \rceil$, greedy tree construction.

Huffman code: in each step, add the two smallest probabilities. Huffman is optimal among character codings.

Entropy as a measure of information. Theoretical lower bound for average codeword length.

$$H(X) = \sum_i p_i \log_2(1/p_i)$$

Arithmetic coding: partition [0,1] according to character probabilities, repeat for subintervals in each step. Not a character encoding, but asymptotically optimal.

We have a source with the following distribution and code table:

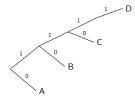
Source symbol	Probability	Codeword
Α	0.4	0
В	0.2	10
C	0.2	110
D	0.2	1111

- (a) Is this a prefix code?
- (b) What is the average codelength?
- (c) How far is the average codelength from the theoretical lower bound of compressibility?
- (d) Is this an optimal character encoding?

Solution.

(a) In order to check if the code is prefix, we can draw the code tree:

Α	0
В	10
С	110
D	1111

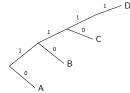


Characters are only in the leaves, so this is a prefix code.

Solution.

(a) In order to check if the code is prefix, we can draw the code tree:

Α	0	
В	10	
С	110	
D	1111	<



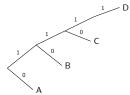
Characters are only in the leaves, so this is a prefix code.

(b)
$$L = \sum_{i=1}^{4} p_i \ell_i = 0.4 \cdot 1 + 0.2 \cdot 2 + 0.2 \cdot 3 + 0.2 \cdot 4 = 2.2.$$

Solution.

(a) In order to check if the code is prefix, we can draw the code tree:

Α	0	
В	10	
С	110	1
D	1111	



Characters are only in the leaves, so this is a prefix code.

- (b) $L = \sum_{i=1}^{4} p_i \ell_i = 0.4 \cdot 1 + 0.2 \cdot 2 + 0.2 \cdot 3 + 0.2 \cdot 4 = 2.2.$
- (c) Theoretical lower bound for average codeword length:

$$H(X) = \sum_{i=1}^{4} p_i \log_2 \left(\frac{1}{p_i}\right) = 0.4 \cdot 1.31 + 3 \cdot 0.2 \cdot 2.322 = 1.922,$$

SO

$$L - H(X) = 0.278$$

(d) Is this an optimal character encoding?

(d) Is this an optimal character encoding?

For the character D, the codeword 111 is sufficient instead of 1111 (and the code is still prefix), so the code is not optimal.

(d) Is this an optimal character encoding?

For the character D, the codeword 111 is sufficient instead of 1111 (and the code is still prefix), so the code is not optimal.

What if we use D \rightarrow 111, is the code optimal then?

(d) Is this an optimal character encoding?

For the character D, the codeword 111 is sufficient instead of 1111 (and the code is still prefix), so the code is not optimal.

What if we use D ightarrow 111, is the code optimal then?

(d) Is this an optimal character encoding?

For the character D, the codeword 111 is sufficient instead of 1111 (and the code is still prefix), so the code is not optimal.

What if we use D ightarrow 111, is the code optimal then?

$$p_1 = 0.4$$

$$p_2 = 0.2$$

$$p_3 = 0.2$$

$$p_4 = 0.2$$

(d) Is this an optimal character encoding?

For the character D, the codeword 111 is sufficient instead of 1111 (and the code is still prefix), so the code is not optimal.

What if we use D ightarrow 111, is the code optimal then?

$$p_1 = 0.4$$
 0.4
 $p_2 = 0.2$ 0.2
 $p_3 = 0.2$ 0.4
 $p_4 = 0.2$

(d) Is this an optimal character encoding?

For the character D, the codeword 111 is sufficient instead of 1111 (and the code is still prefix), so the code is not optimal.

What if we use D ightarrow 111, is the code optimal then?

$$p_1 = 0.4$$
 0.4 0.4 $p_2 = 0.2$ 0.6 $p_3 = 0.2$ 0.4 $p_4 = 0.2$

(d) Is this an optimal character encoding?

For the character D, the codeword 111 is sufficient instead of 1111 (and the code is still prefix), so the code is not optimal.

What if we use D ightarrow 111, is the code optimal then?

$$\rho_1 = 0.4$$
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 0.4
 $0.$

(d) Is this an optimal character encoding?

For the character D, the codeword 111 is sufficient instead of 1111 (and the code is still prefix), so the code is not optimal.

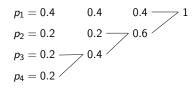
What if we use D \rightarrow 111, is the code optimal then?

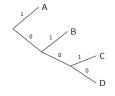
$$p_1 = 0.4$$
 0.4 0.4 1
 $p_2 = 0.2$ 0.2 0.6
 $p_3 = 0.2$ 0.4
 $p_4 = 0.2$

(d) Is this an optimal character encoding?

For the character D, the codeword 111 is sufficient instead of 1111 (and the code is still prefix), so the code is not optimal.

What if we use D \rightarrow 111, is the code optimal then?





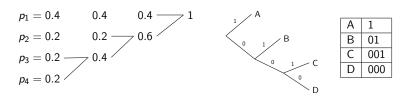
Α	1
В	01
С	001
D	000

(d) Is this an optimal character encoding?

For the character D, the codeword 111 is sufficient instead of 1111 (and the code is still prefix), so the code is not optimal.

What if we use D ightarrow 111, is the code optimal then?

We can check by comparing it to Huffman:



Same codeword lengths \to the code is equivalent to Huffman, both are optimal, and the average codeword length is

$$L = 0.4 \cdot 1 + 0.2 \cdot 2 + 0.2 \cdot 3 + 0.2 \cdot 3 = 2.0.$$



(d) During Huffman, in the second column we picked the bottom 0.4.

$$\rho_1 = 0.4$$
 0.4
 0.4
 0.4
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 0.6
 $0.$

(d) During Huffman, in the second column we picked the bottom 0.4.

$$p_1 = 0.4$$
 0.4 0.4 1
 $p_2 = 0.2$ 0.2 0.6
 $p_3 = 0.2$ 0.4

What if we pick the other 0.4?

$$p_1 = 0.4$$
 0.4 0.6 1
 $p_2 = 0.2$ 0.2
 $p_3 = 0.2$ 0.4 0.4
 $p_4 = 0.2$

(d) During Huffman, in the second column we picked the bottom 0.4.

$$p_1 = 0.4$$
 0.4 0.4 1
 $p_2 = 0.2$ 0.2 0.6
 $p_3 = 0.2$ 0.4

What if we pick the other 0.4?

$$p_1 = 0.4$$
 0.4 0.6 1
 $p_2 = 0.2$ 0.2 0.4 0.4 $p_3 = 0.2$ 0.4 0.4

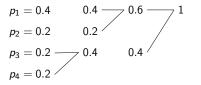


Α	11
В	10
C	01
D	00

(d) During Huffman, in the second column we picked the bottom 0.4.

$$p_1 = 0.4$$
 0.4 0.4 1
 $p_2 = 0.2$ 0.2 0.6
 $p_3 = 0.2$ 0.4

What if we pick the other 0.4?





11
10
01
00

L=2, same as for the previous choice. This coding is also an optimal Huffman code, just using a different tree.

Encode the following distribution using Fixed length coding, Shannon–Fano coding and Huffman coding. Compute the average codeword length for each coding. Compute the theoretical lower bound.

$$p_1 = 0.49$$
 $p_2 = 0.14$
 $p_3 = 0.14$
 $p_4 = 0.07$
 $p_5 = 0.07$
 $p_6 = 0.04$
 $p_7 = 0.02$
 $p_8 = 0.02$
 $p_9 = 0.01$

Solution. For fixed length encoding, the required codeword length is $\lceil \log_2 9 \rceil = 4,$ and

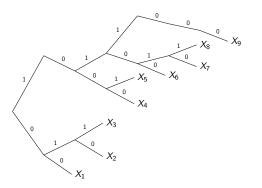
Symbol	Codeword
X_1	0000
X_2	0001
<i>X</i> ₃	0010
X_4	0011
<i>X</i> ₅	0100
X_6	0101
X_7	0110
<i>X</i> ₈	0111
X_9	1000

Average codeword length is L = 4.

For the Shannon–Fano coding, the codeword lengths are $\ell_i = \lceil \log_2 1/p_i \rceil$, so

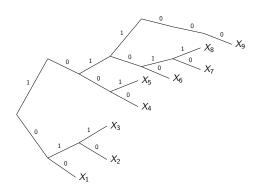
$$\begin{split} \ell_1 &= \lceil \log_2 1/0.49 \rceil = \lceil 1.029 \rceil = 2, \\ \ell_2 &= \lceil \log_2 1/0.14 \rceil = \lceil 2.836 \rceil = 3, \\ \ell_3 &= \lceil \log_2 1/0.04 \rceil = \lceil 2.836 \rceil = 3, \\ \ell_4 &= \lceil \log_2 1/0.07 \rceil = \lceil 3.836 \rceil = 4, \\ \ell_5 &= \lceil \log_2 1/0.07 \rceil = \lceil 3.836 \rceil = 4, \\ \ell_6 &= \lceil \log_2 1/0.04 \rceil = \lceil 4.644 \rceil = 5, \\ \ell_7 &= \lceil \log_2 1/0.02 \rceil = \lceil 5.644 \rceil = 6, \\ \ell_8 &= \lceil \log_2 1/0.02 \rceil = \lceil 5.644 \rceil = 6, \\ \ell_9 &= \lceil \log_2 1/0.01 \rceil = \lceil 6.644 \rceil = 7. \end{split}$$

Code tree and codeword table for Shannon-Fano:



Symbol	Codeword
X_1	00
X_2	010
<i>X</i> ₃	011
<i>X</i> ₄	1000
X_5	1001
X_6	10100
<i>X</i> ₇	101010
<i>X</i> ₈	101011
X_9	1011000

Code tree and codeword table for Shannon-Fano:



Symbol	Codeword
X_1	00
X_2	010
<i>X</i> ₃	011
X_4	1000
X_5	1001
X_6	10100
X_7	101010
<i>X</i> ₈	101011
X_9	1011000

The average codeword length is

$$L_{S-F} = 0.49 \cdot 2 + 2 \times 0.14 \cdot 3 + 2 \times 0.07 \cdot 4 +$$

 $0.04 \cdot 5 + 2 \times 0.02 \cdot 6 + 0.01 \cdot 7 = 2.89.$

- $p_1 = 0.49$
- $p_2 = 0.14$
- $p_3 = 0.14$
- $p_4 = 0.07$
- $p_5 = 0.07$
- $p_6 = 0.04$
- $p_7 = 0.02$
- $p_8 = 0.02$
- $p_9 = 0.01$

$$p_1 = 0.49$$
 0.49
 $p_2 = 0.14$ 0.14
 $p_3 = 0.14$ 0.07
 $p_4 = 0.07$ 0.07
 $p_5 = 0.07$ 0.07
 $p_6 = 0.04$ 0.04
 $p_7 = 0.02$ 0.02
 $p_8 = 0.02$ 0.03

$$p_1 = 0.49$$
 0.49 0.49
 $p_2 = 0.14$ 0.14 0.14
 $p_3 = 0.14$ 0.14 0.14
 $p_4 = 0.07$ 0.07 0.07
 $p_5 = 0.07$ 0.07 0.07
 $p_6 = 0.04$ 0.04 0.04
 $p_7 = 0.02$ 0.02 0.05
 $p_8 = 0.02$ 0.03

$$p_1 = 0.49$$
 0.49 0.49 0.49 $p_2 = 0.14$ 0.14 0.14 0.14 0.14 0.14 $p_3 = 0.14$ 0.14 0.14 0.14 0.14 $p_4 = 0.07$ 0.07 0.07 0.07 0.07 $p_5 = 0.07$ 0.07 0.07 0.07 0.09 $p_7 = 0.02$ 0.02 0.03 0.03 $p_9 = 0.01$

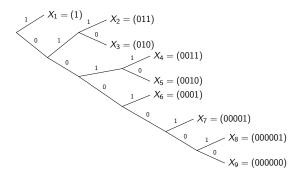
$$\rho_1 = 0.49$$
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.40
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14
 0.14

$$p_1 = 0.49$$
 0.49 0.49 0.49 0.49 0.49 0.49 0.49 $p_2 = 0.14$ 0.14 0.14 0.14 0.14 0.14 0.14 0.28 $p_3 = 0.14$ 0.14 0.14 0.14 0.14 0.14 0.14 0.23 $p_4 = 0.07$ 0.07 0.07 0.07 0.07 0.09 0.09 $p_5 = 0.04$ 0.04 0.04 0.09 $p_7 = 0.02$ 0.03 $p_9 = 0.01$

$$\begin{aligned}
 \rho_1 &= 0.49 & 0.49 & 0.49 & 0.49 & 0.49 & 0.49 & 0.49 \\
 \rho_2 &= 0.14 & 0.14 & 0.14 & 0.14 & 0.14 & 0.14 & 0.28 & 0.51 \\
 \rho_3 &= 0.14 & 0.14 & 0.14 & 0.14 & 0.14 & 0.14 & 0.23 & 0.23 \\
 \rho_4 &= 0.07 & 0.07 & 0.07 & 0.07 & 0.04 & 0.23 \\
 \rho_5 &= 0.07 & 0.07 & 0.07 & 0.09 & 0.09 & 0.02 & 0.02 & 0.05 \\
 \rho_6 &= 0.04 & 0.04 & 0.04 & 0.09 & 0.09 & 0.05 & 0.03 & 0.05 \\
 \rho_8 &= 0.02 & 0.03 & 0.03 & 0.03 & 0.09 & 0.01 & 0.01 & 0.01 & 0.01 \\
 \rho_9 &= 0.01 & 0.03 & 0$$

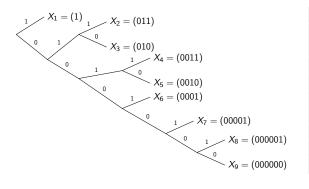
$$\rho_1 = 0.49$$
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49
 0.49

Then the code tree and codeword table can be obtained:



Symbol	Codeword
X_1	1
X_2	011
<i>X</i> ₃	010
<i>X</i> ₄	0011
X_5	0010
X_6	0001
<i>X</i> ₇	00001
<i>X</i> ₈	000001
<i>X</i> ₉	000000

Then the code tree and codeword table can be obtained:



Symbol	Codeword		
X_1	1		
X_2	011		
<i>X</i> ₃	010		
X_4	0011		
X_5	0010		
X_6	0001		
<i>X</i> ₇	00001		
<i>X</i> ₈	000001		
X_9	000000		

$$L_{Huff} = 0.49 \cdot 1 + 0.14 \cdot 3 + 0.14 \cdot 3 + 0.07 \cdot 4 + 0.07 \cdot 4 + 0.04 \cdot 4 + 0.02 \cdot 5 + 0.02 \cdot 6 + 0.01 \cdot 6 = 2.33$$

Entropy of the source:

$$H(X) = \sum_{i=1}^{9} p_i \log_2 \left(\frac{1}{p_i}\right) = 2.314.$$

Entropy of the source:

$$H(X) = \sum_{i=1}^{9} p_i \log_2 \left(\frac{1}{p_i}\right) = 2.314.$$

Comparison with the average codeword length for the various codings:

$$L_{fixed} = 4$$

 $L_{S-F} = 2.89$
 $L_{Huff} = 2.33$

A source has alphabet $\{A,B,C,D\}$, with distribution

$$P(A) = 0.4$$
, $P(B) = 0.3$, $P(C) = 0.2$, $P(D) = 0.1$.

For each of Fixed length coding, Huffman coding and Arithmetic coding, what is the length of the codeword corresponding to the message AAAAA?

And for the messages BBBBB, CCCCC and DDDDD?

A source has alphabet $\{A,B,C,D\}$, with distribution

$$P(A) = 0.4$$
, $P(B) = 0.3$, $P(C) = 0.2$, $P(D) = 0.1$.

For each of Fixed length coding, Huffman coding and Arithmetic coding, what is the length of the codeword corresponding to the message AAAAA?

And for the messages BBBBB, CCCCC and DDDDD?

Solution. For Fixed length encoding, all characters are 2 bits, so the length of each of the four messages is 10 bits.



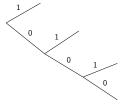
- $p_1 = 0.4$
- $p_2 = 0.3$
- $p_3 = 0.2$
- $p_4 = 0.1$

$$p_1 = 0.4$$
 0.4
 $p_2 = 0.3$ 0.3
 $p_3 = 0.2$ 0.3
 $p_4 = 0.1$

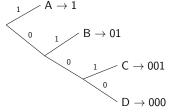
$$p_1 = 0.4$$
 0.4 0.4 0.4 $p_2 = 0.3$ 0.3 0.6 $p_3 = 0.2$ 0.3 0.3

$$p_1 = 0.4$$
 0.4 0.4 1
 $p_2 = 0.3$ 0.3 0.6
 $p_3 = 0.2$ 0.3
 $p_4 = 0.1$

$$p_1 = 0.4$$
 0.4 0.4 1
 $p_2 = 0.3$ 0.3 0.6
 $p_3 = 0.2$ 0.3



$$p_1 = 0.4$$
 0.4 0.4 1
 $p_2 = 0.3$ 0.3 0.6
 $p_3 = 0.2$ 0.3



For Huffman coding,

The codewords are:

- ► AAAAA → 11111, length 5
- ▶ BBBBB \rightarrow 0101010101, length 10
- ightharpoonup CCCCC ightharpoonup 001001001001001, length 15
- ▶ DDDDD → 00000000000000, length 15

For Arithmetic coding, the length of the messages:

$$\begin{array}{lll} \mathsf{AAAAA} & \to & \left\lceil -\log_2\left(0.4^5\right) \right\rceil + 1 = 7 \\ \\ \mathsf{BBBBB} & \to & \left\lceil -\log_2\left(0.3^5\right) \right\rceil + 1 = 10 \\ \\ \mathsf{CCCCC} & \to & \left\lceil -\log_2\left(0.2^5\right) \right\rceil + 1 = 13 \\ \\ \mathsf{DDDDD} & \to & \left\lceil -\log_2\left(0.1^5\right) \right\rceil + 1 = 18 \end{array}$$

Using the Fixed length encoding, Huffman and Arithmetic coding from the previous problem, compute the codeword length for a message containing 40 A's, 30 B's, 20 C's and 10 D's total.

Using the Fixed length encoding, Huffman and Arithmetic coding from the previous problem, compute the codeword length for a message containing 40 A's, 30 B's, 20 C's and 10 D's total.

Solution. For Fixed length encoding, the codeword length is

$$100 \times 2 = 200$$
 bits.

Using the Fixed length encoding, Huffman and Arithmetic coding from the previous problem, compute the codeword length for a message containing 40 A's, 30 B's, 20 C's and 10 D's total.

Solution. For Fixed length encoding, the codeword length is

$$100 \times 2 = 200$$
 bits.

Using Huffman, it is

$$40 \times 1 + 30 \times 2 + 20 \times 3 + 10 \times 3 = 190$$
 bits.



Using the Fixed length encoding, Huffman and Arithmetic coding from the previous problem, compute the codeword length for a message containing 40 A's, 30 B's, 20 C's and 10 D's total.

Solution. For Fixed length encoding, the codeword length is

$$100 \times 2 = 200$$
 bits.

Using Huffman, it is

$$40 \times 1 + 30 \times 2 + 20 \times 3 + 10 \times 3 = 190$$
 bits.

With Arithmetic coding, it is

$$\left[-\log_2(0.4^{40}\cdot 0.3^{30}\cdot 0.2^{20}\cdot 0.1^{10})\right]+1=186 \text{ bits.}$$



We have a source with 7 characters and probabilities

$$p_1 = 0.5$$
, $p_2 = p_3 = 0.125$, $p_4 = p_5 = p_6 = p_7 = 0.0625$.

Compute the codeword lengths for this source for both Shannon–Fano coding and Huffman coding. Also compute the entropy of the source.

We have a source with 7 characters and probabilities

$$p_1 = 0.5$$
, $p_2 = p_3 = 0.125$, $p_4 = p_5 = p_6 = p_7 = 0.0625$.

Compute the codeword lengths for this source for both Shannon–Fano coding and Huffman coding. Also compute the entropy of the source.

Solution. For the Shannon–Fano coding, the codeword lengths are $\ell_i = \lceil \log_2 1/p_i \rceil$, that is,

$$\ell_1 = 1, \quad \ell_2 = \ell_3 = 3, \quad \ell_4 = \ell_5 = \ell_6 = \ell_7 = 4,$$

and the average codeword length is

$$L_{S-F} = 0.5 \times 1 + 2 \times 0.125 \times 3 + 4 \times 0.0625 \times 4 = 2.25.$$



For Huffman coding,

0.5

0.125

0.125

0.0625

0.0625

0.0625

0.0625

0.5	1
0.125	0.125
0.125	0.125
0.0625	0.0625
0.0625	0.0625
0.0625	0.125
0.0625	

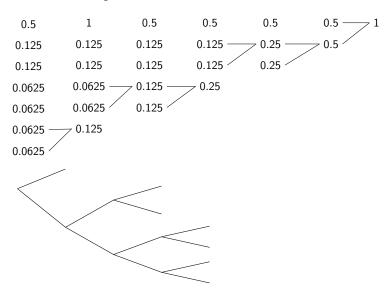
0.5	1	0.5
0.125	0.125	0.125
0.125	0.125	0.125
0.0625	0.0625	0.125
0.0625	0.0625	0.125
0.0625	0.125	
0.0625		

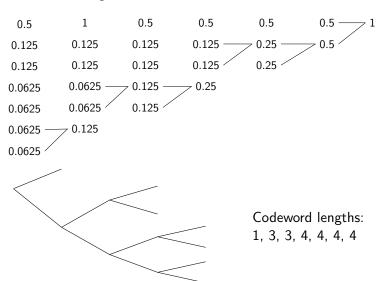
```
0.5
                         0.5
                                     0.5
            0.125
                        0.125
                                    0.125
0.125
0.125
            0.125
                        0.125
                                    0.125
            0.0625
                        0.125
                                    0.25
0.0625
                        0.125
            0.0625
0.0625
            0.125
0.0625 -
0.0625
```

```
0.5
                         0.5
                                     0.5
                                                 0.5
0.125
            0.125
                        0.125
                                    0.125
                                                0.25
                        0.125
                                    0.125
                                                0.25
0.125
            0.125
                                    0.25
           0.0625
                        0.125
0.0625
                        0.125
           0.0625
0.0625
            0.125
0.0625
0.0625
```

```
0.5
                         0.5
                                     0.5
                                                 0.5
                                                            0.5
0.125
            0.125
                        0.125
                                    0.125
                                                0.25
                                                             0.5
                        0.125
                                    0.125
                                                0.25
0.125
            0.125
                                    0.25
           0.0625
                        0.125
0.0625
                        0.125
           0.0625
0.0625
            0.125
0.0625
0.0625
```

```
0.5
                         0.5
                                     0.5
                                                 0.5
                                    0.125
                                                0.25
0.125
            0.125
                        0.125
                        0.125
                                    0.125
                                                0.25
0.125
            0.125
                                    0.25
           0.0625
                        0.125
0.0625
                        0.125
           0.0625
0.0625
            0.125
0.0625
0.0625
```





For Huffman coding, the average codeword length is

$$L_{Huff} = 0.5 \times 1 + 2 \times 0.125 \times 3 + 4 \times 0.0625 \times 4 = 2.25.$$

For Huffman coding, the average codeword length is

$$L_{Huff} = 0.5 \times 1 + 2 \times 0.125 \times 3 + 4 \times 0.0625 \times 4 = 2.25.$$

The entropy of the source is

$$H(X) = -(0.5 \times \log_2(0.5) + 2 \times 0.125 \cdot \log_2(0.125) + 4 \times 0.0625 \cdot \log_2(0.0625)) = 2.25.$$

For Huffman coding, the average codeword length is

$$L_{Huff} = 0.5 \times 1 + 2 \times 0.125 \times 3 + 4 \times 0.0625 \times 4 = 2.25.$$

The entropy of the source is

$$H(X) = -(0.5 \times \log_2(0.5) + 2 \times 0.125 \cdot \log_2(0.125) + 4 \times 0.0625 \cdot \log_2(0.0625)) = 2.25.$$

Both Shannon–Fano and Huffman reach the theoretical lower bound for this source.

For Huffman coding, the average codeword length is

$$L_{\text{Huff}} = 0.5 \times 1 + 2 \times 0.125 \times 3 + 4 \times 0.0625 \times 4 = 2.25.$$

The entropy of the source is

$$H(X) = -(0.5 \times \log_2(0.5) + 2 \times 0.125 \cdot \log_2(0.125) + 4 \times 0.0625 \cdot \log_2(0.0625)) = 2.25.$$

Both Shannon–Fano and Huffman reach the theoretical lower bound for this source.

Actually, Huffman and Shannon–Fano both reach the entropy for the same sources in general: when all probabilities are negative integer powers of 2.

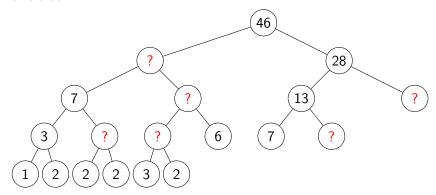
Adaptive Huffman code

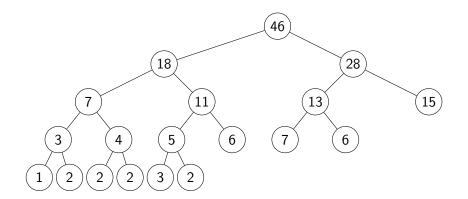
Adaptive Huffman code: the code tree is based on the number of occurrences of each character seen so far. Internal nodes get the sum of their children.

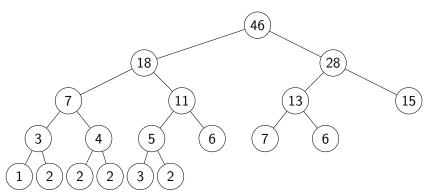
- 1. Initialize the code tree.
- 2. Code the next character according to the current state of the code tree.
- 3. Add the character to the tree.
- 4. Check the sibling property, and restore it if necessary by exchanging two nodes with their entire subtrees.
- 5. Go to next character and repeat from step 2.

Sibling property: listing the nodes from bottom level to top level, going from left to right within each level, the nodes are increasing, except possibly for sibling pairs.

Determine the missing values. Does the sibling property hold for this tree?



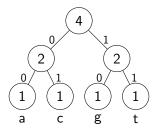




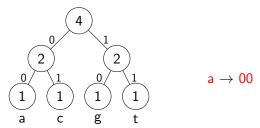
Starting from the bottom left, going left to right first, then up one level and repeat, the sequence (1, 2), (2, 2), (3, 2), (3, 4), (5, 6), (7, 6), (7, 11), (13, 15), (18, 28), 46 is decreasing only for two sibling pairs, so the sibling property holds for this tree, this is a valid Huffman-tree.

Over the alphabet $\{a,c,g,t\}$, apply adaptive Huffman coding to the message "aagacaa".

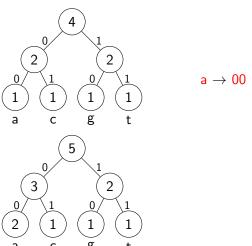
Over the alphabet $\{a,c,g,t\}$, apply adaptive Huffman coding to the message "aagacaa".



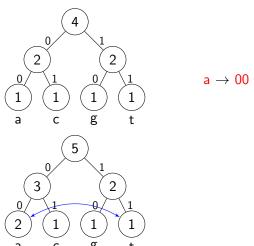
Over the alphabet $\{a,c,g,t\}$, apply adaptive Huffman coding to the message "aagacaa".



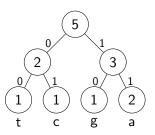
Over the alphabet $\{a,c,g,t\}$, apply adaptive Huffman coding to the message "aagacaa".



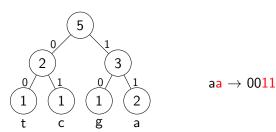
Over the alphabet $\{a,c,g,t\}$, apply adaptive Huffman coding to the message "aagacaa".



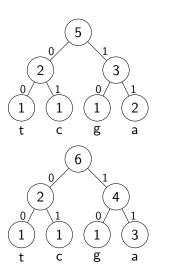
aagacaa



aagacaa

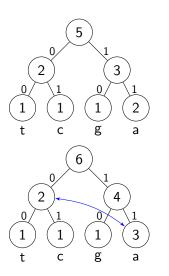


aagacaa



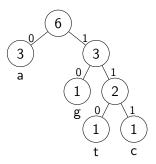
 $aa \rightarrow 0011$

aagacaa

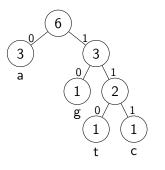


 $a\textbf{a} \rightarrow 0011$

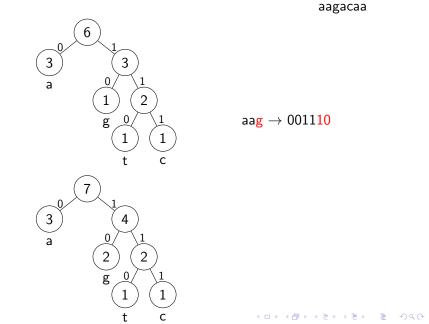
aagacaa

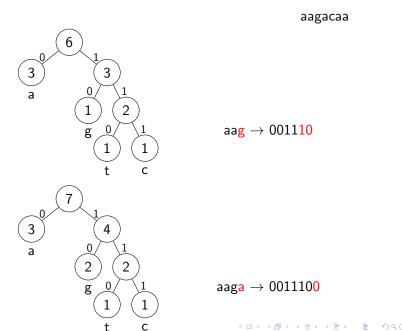




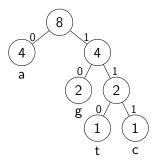


$$\mathsf{aag} \to \mathsf{001110}$$

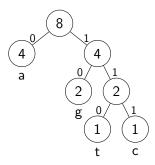




aagacaa

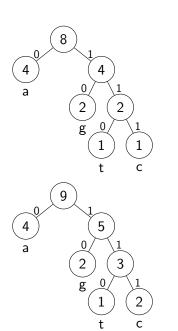


aagacaa

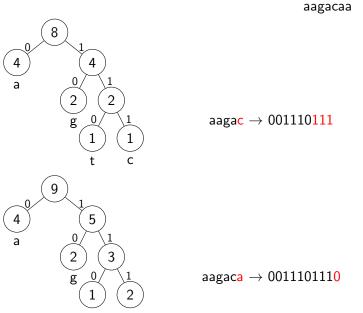


 $\mathsf{aagac} \to 001110\textcolor{red}{111}$

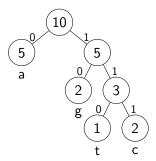




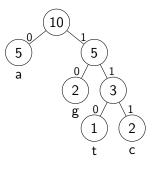
 $\mathsf{aagac} \to \mathsf{001110}\textcolor{red}{\textbf{111}}$



aagacaa

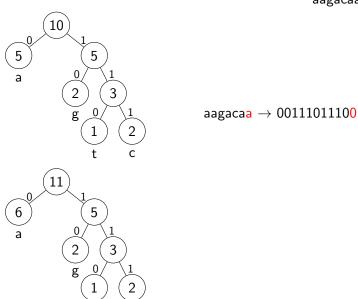


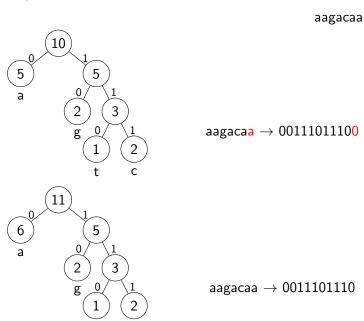
aagacaa



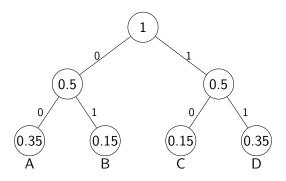
 $\mathsf{aagaca} \to \mathsf{00111011100}$

aagacaa

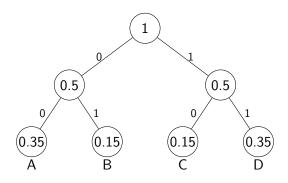




Can we obtain the following code tree as the result of a Huffman algorithm for code tree construction?

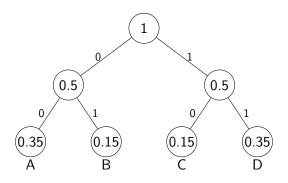


Can we obtain the following code tree as the result of a Huffman algorithm for code tree construction?



Solution 1. No, the sibling property does not hold.

Can we obtain the following code tree as the result of a Huffman algorithm for code tree construction?

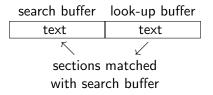


Solution 1. No, the sibling property does not hold.

Solution 2. For the Huffman tree construction, in the first step of the algorithm, the two smallest probabilities are the two 0.15's, so they would have to be matched to each other, not to the 0.35's.

Reminder: Dictionary coders

LZ77: search ahead for a section of text seen recently.



Output is a sequence of (p, ℓ, c) records (position of match, length of match, next character).

LZ78: parse the text into sections 1 character longer than seen before. The output is a sequence of records (i, c) (address of old section, new character).

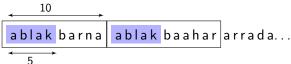
Compute the next two records of the LZ77 algorithm, starting from the following position.

bablakbarna ablakbaahar arrada...

Compute the next two records of the LZ77 algorithm, starting from the following position.

Solution.

first step:



Compute the next two records of the LZ77 algorithm, starting from the following position.

Solution.

first step:



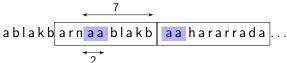
Compute the next two records of the LZ77 algorithm, starting from the following position.

Solution.

first step:



next step:



Compute the next two records of the LZ77 algorithm, starting from the following position.

Solution.

first step:



next step:

An LZ77 coder has parameters $h_s = h_\ell = 4$. Decode the following sequence of records: (0,0,C), (0,0,A), (2,1,D), (3,2,B), (2,5,A).

An LZ77 coder has parameters $h_s = h_\ell = 4$. Decode the following sequence of records: (0,0,C), (0,0,A), (2,1,D), (3,2,B), (2,5,A).

Solution. (0,0,X)-type records code a single character (X), so the first two records are decoded to CA.

An LZ77 coder has parameters $h_s = h_\ell = 4$. Decode the following sequence of records: (0,0,C), (0,0,A), (2,1,D), (3,2,B), (2,5,A).

Solution. (0,0,X)-type records code a single character (X), so the first two records are decoded to CA.

The numbers in the record (2,1,D) are decoded as a copy command: move back 2 positions and copy 1 character (C). Then the D is added to the end \rightarrow CACD.

An LZ77 coder has parameters $h_s = h_\ell = 4$. Decode the following sequence of records: (0,0,C), (0,0,A), (2,1,D), (3,2,B), (2,5,A).

Solution. (0,0,X)-type records code a single character (X), so the first two records are decoded to CA.

The numbers in the record (2,1,D) are decoded as a copy command: move back 2 positions and copy 1 character (C). Then the D is added to the end \rightarrow CACD.

Similarly, (3,2,B) results in CACDACB.

An LZ77 coder has parameters $h_s = h_\ell = 4$. Decode the following sequence of records: (0,0,C), (0,0,A), (2,1,D), (3,2,B), (2,5,A).

Solution. (0,0,X)-type records code a single character (X), so the first two records are decoded to CA.

The numbers in the record (2,1,D) are decoded as a copy command: move back 2 positions and copy 1 character (C). Then the D is added to the end \rightarrow CACD.

Similarly, (3,2,B) results in CACDACB.

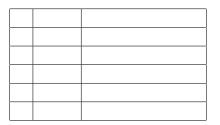
For the record (2,5,A), we execute a similar copy command, but we only have 2 characters (CB) to copy; in this case, we keep copying those 2 characters cyclically until we get 5 characters (CBCBC), then add an A to the end \rightarrow CACDACBCBCBCA.

Using LZ78, encode the message AABABBABBABB. Also convert the dictionary to binary.

Using LZ78, encode the message AABABBABBABB. Also convert the dictionary to binary.

Solution.

AABABBBABBABB



Using LZ78, encode the message AABABBABBABB. Also convert the dictionary to binary.

Solution.

 $\mathsf{A}|\mathsf{A}\,\mathsf{B}\,\mathsf{A}\,\mathsf{B}\,\mathsf{B}\,\mathsf{B}\,\mathsf{B}\,\mathsf{A}\,\mathsf{B}\,\mathsf{B}\,\mathsf{A}\,\mathsf{B}\,\mathsf{B}\,\mathsf{A}$

1	(<mark>0,A</mark>)	1

Using LZ78, encode the message AABABBABBABB. Also convert the dictionary to binary.

Solution.

$A|A\,B|A\,B\,B\,B\,A\,B\,B\,A\,B\,B$

1	(0,A)	1
2	(1,B)	111

Using LZ78, encode the message AABABBABBABB. Also convert the dictionary to binary.

Solution.

$A|A\,B|A\,B\,B|B\,A\,B\,B\,A\,B\,B$

1	(0,A)	1
2	(1,B)	111
3	(2,B)	111101

Using LZ78, encode the message AABABBABBABB. Also convert the dictionary to binary.

Solution.

$A|A\,B|A\,B\,B|B|A\,B\,B\,A\,B\,B$

1	(0,A)	1
2	(1,B)	111
3	(2,B)	111101
4	(0,B)	11110001

Using LZ78, encode the message AABABBABBABB. Also convert the dictionary to binary.

Solution.

$A|A\ B|A\ B\ B|B|A\ B\ B\ A|B\ B$

1	(0,A)	1
2	(1,B)	111
3	(2,B)	111101
4	(0,B)	11110001
5	(3,A)	1111000 <mark>0110</mark>

Using LZ78, encode the message AABABBABBABB. Also convert the dictionary to binary.

Solution.

$A|A\ B|A\ B\ B|B|A\ B\ B\ A|B\ B$

1	(0,A)	1
2	(1,B)	111
3	(2,B)	111101
4	(0,B)	11110001
5	(3,A)	1111000 <mark>0110</mark>
6	(4,B)	111100001101001