

ELEN060-2 - Information and coding theory

Project 2 - Source coding, data compression and channel coding

March 2022

The goal of this second project is to apply some of the principles seen in the lectures about source coding, data compression and channel coding. We ask you to write a brief report (pdf format) collecting your answers to the different questions. All codes must be written in Python inside the Jupyter Notebook provided with this assignment, no other code file will be accepted. Note that you can not change the content of locked cells or import any extra Python library than the ones provided.

The assignment must be carried out by groups of two students and the report and the notebook should be submitted on Gradescope (<https://www.gradescope.com/>) before May 4 23:59. Note that attention will be paid to how you present your results and your analyses. By submitting the project, each member of a group shares the responsibility for what has been submitted (e.g., in case of plagiarism). From a practical point of view, every student should have registered on the platform before the deadline. Group, archive and report should be named by the concatenation of your student ID (sXXXXXX) (e.g., s000007s123456.pdf and s000007s123456.ipynb).

Implementation

1. Implement a function that returns a *binary* Huffman code for a given probability distribution. Give the main steps of your implementation. Verify your code on Exercise 7 of the second exercise session (TP2), and report the output of your code for this example. Explain how to extend your function to generate a Huffman code of *any (output) alphabet size*.
2. Given a sequence of symbols, implement a function that returns a dictionary and the encoded sequence using the *on-line Lempel-Ziv algorithm* (see [State of the art in data compression](#), slide 50/53). Reproduce and report the example given in the course.
3. Compare (without implementing the basic version) the two versions of the Lempel-Ziv algorithm seen in the theoretical course (i.e., the basic and the on-line versions). Discuss what are the (practical) advantages and drawbacks of each version.
4. The LZ77 algorithm is another dictionary algorithm. It consists in replacing repeated sequences of symbols with references to a previous occurrence of the same sequence of symbols. Typically, the encoder considers a sliding window search buffer of size *window_size* and a look-ahead buffer (you can constrain the size of the look-ahead buffer if you wish, but pay attention to choose a relevant size and to specify it in your report), and searches for past occurrences of the beginning of the

look-ahead buffer (i.e, the prefix) within the search buffer. Each codeword is made of three elements: the offset (i.e, distance) of the longest prefix in the search buffer, the length of the prefix, and the symbol following the prefix. Note that the offset and the length are in practice rewritten in binary.

Implement a function that returns the encoded sequence using the LZ77 algorithm as described by the algorithm below given an input string and a sliding window size. Reproduce the example given in Figure 2 with *window_size*=7.

Algorithm 1: LZ77 compression algorithm sliding window

Result: Write here the result

A sliding window size l ;

An input string;

while *input is not empty* **do**

 prefix := longest prefix of input that begins in window;

if *prefix exists in window* **then**

d := distance to the start of the prefix;

p := length of prefix;

c := char following the prefix in input;

else

d := 0;

p := 0;

c := first symbol of input;

end

 append (d, p, c) to encoded input;

 shift the sliding window by $p + 1$ symbols (*i.e.*, discard $p + 1$ symbols from the beginning of window and add the $p + 1$ first symbols of the input at the end of the window).

end

7	6	5	4	3	2	1										output
						a	a	b	r	a	c	a	ada...			(0,0,a)
					a	b	b	r	a	c	a	d	dab...			(0,0,b)
				a	b	r	a	c	a	d	a		abr...			(0,0,r)
		a	b	r	a	c	a	d	a	b	r	a	bra...			(3,1,c)
							a	d	a	b	r	a	ad...			(2,1,d)
a	b	r	a	c	a	d	a	b	r	a	d					(7,4,d)
a	d	a	b	r	a	d										
sliding window							look-ahead buffer									

Figure 2 - Example of the encoding of the string *abracadabrad* using the LZ77 algorithm. Taken from this source: <http://jens.jm-s.de/comp/LZ77-JensMueller.pdf>

Source coding and reversible (lossless) data compression

Before the invention of the telephones to communicate over long distances, Samuel F. B. Morse, along with Joseph Henry and Alfred Vail, invented the telegraph. This machine used electricity to transfer a message through the wires. This message was a sequence of short

and long signals, known as *dots* and *dashes* in the Morse code. This character-encoding scheme assigns to each letter a unique combination of *dots* and *dashes* signals. Later on, Morse code had an extensive usage during World War II. Table 1 summarizes the encoding system.

For the purpose of this project we consider a four symbol morse code by adding a word separator by a *dash symbol* “/” and letter separator by an *underscore* “_”. We also don’t consider punctuation marks in the original text (dots, question marks, etc.) or digits. With this assignment, you are given a *morse.txt* file containing a text encoded in Morse code (4 symbols) and the original text with 27 symboles (26 letters and word space) in *text.txt*.

International Morse alphabet													
A	B	C	D	E	F	G	H	I	J	K	L	M	letter space
.-	-...	-.-.	-..	.	..-.	--.-.-.	-.-	.-..	--	_
N	O	P	Q	R	S	T	U	V	W	X	Y	Z	word space
-.	---	.-.-.	---.	.-.	...	-	..-	...-	.-.-	-.-.-	-.--	---.	/

Table 1 - International Morse Code for the 26 alphabet letters + the two separator symbols we introduce.

- Estimate the marginal probability distribution of all symbols (.,-,_,/) from the given Morse text, and determine the corresponding binary Huffman code and the encoded Morse text. Give the total length of the encoded Morse text and the compression rate.
- Give the expected average length for your Huffman code. Compare this value with (a) the empirical average length, and (b) theoretical bound(s). Justify.
- Plot the evolution of the empirical average length of the encoded morse using your Huffman code for increasing input text lengths.
- Encode the morse code using the *on-line Lempel-Ziv algorithm*. Give the total length of the encoded code and the compression rate.
- Encode the text using the *LZ77 algorithm* with *window_size=7*. Give the total length of the encoded text and the compression rate.
- Famous data compression algorithms combine the LZ77 algorithm and the Huffman algorithm. Explain two ways of combining those algorithms and discuss the interest of the possible combinations.
- Encode the morse using one of the combinations of LZ77 and Huffman algorithms you proposed in the previous question. Give the total length of the encoded morse text and the compression rate.
- Report the total lengths and compression rates using (a) LZ77 and (b) the combination of LZ77 and Huffman, to encode the morse code for different values of the sliding window size (use sliding window sizes from 1 to 11000 with a step of 1000). Compare your result with the total length and compression rate obtained using the on-line Lempel-Ziv algorithm. Discuss your results.
- It is typically assumed that repetitions occur at long distances in a text (for instance the name of a character in a book). Based on your results in the previous question(s), discuss what could be the best data compression algorithm(s) and/or how to adapt the algorithms used in this project.

14. Instead of encoding the Morse code, encode directly the original text (27 symbols) with the binary Huffman algorithm. Give the code for each symbol, the average expected length, the experimental length of the encoded text and the compression rate.
15. Compare the values found at the previous question with the ones found in Question 5. In particular, is it better to first encode the text with Morse code before the Huffman encoding or to directly encode the text with Huffman? Discuss.

Channel coding

Let us consider a grayscale image that is sent through a noisy channel. Let us take a .png file *image.png* as an image signal. Its quantisation is such that possible values are between 0 and 255, and its number of pixels is *width* \times *height*. The channel is a binary symmetric channel with a probability of error equal to 0.01. In order to send the image signal through the channel, the signal is first encoded in a binary alphabet and then each binary symbol is sent through the channel. The pixel values are read from left to right and from top to bottom as for conventional text reading in English.



16. Implement a function to read the image and display the original image. Remember that the image should be read in grayscale (one channel per pixel).
17. Encode the image signal using a fixed-length binary code. What is the appropriate number of bits? Justify.
18. Simulate the channel effect on the binary image signal. Then decode the image signal and display the decoded image. What do you notice?
19. Instead of sending directly through the channel the binary image signal, you will first introduce some redundancy. To do that, implement a function that returns the Hamming (7,4) code for a given sequence of binary symbols. Then, using your function, encode the binary image signal (from question 16).
20. Simulate the channel effect on the binary image signal with redundancy. Then decode the binary image signal. Display the decoded image signal. What do you notice? Explain your decoding procedure.
21. How would you proceed to reduce the loss of information and/or to improve the communication rate? Justify