



Compte-Rendu

Student name: **TRAN Thanh Duy**

Course : **Codage de l'information**

Teacher : Jean-Christophe Burie



OBJECTIF

Ce TP sert à implémenter l'algorithme du codage LZ78 (Lempel-Ziv) en 2 étapes :

- Encodage et Décodage

Pour l'étape du Encodage, le résultat obtenu est au forme dictionnaire (key-value). Néanmoins, l'étape va retourner : dictionnaire et l'information original (l'information avant encodé).

Les test-cases sont implémenter en deux formes :

1. Des chaines de caractère simples (ex : wabba_wabba_wabba_woo_woo)
2. Un fichier text

La performance de l'algorithme va être évaluer par la ratio de compresse. Cette ration est calculée par la formule (ratio = $\text{size_before_encoded} / \text{size_after_encoded}$)

ALGORITHME DE LZ78 – PHASE D'ENCODAGE

LZ78 – ENCODING PHASE

Tout d'abord, je vais déclarer 2 variables globales servant à sauvegarder les résultats :

Firstly, I will declare 2 global variables aiming to store results

```
comp_dict = []  
decomp_dict = []
```

Comme j'ai déjà exprimé, la phase d'encodage va recevoir un seul paramètre étant la chaîne de caractère (input) et va retourner le résultat sous forme d'un dictionnaire

As mentioned, encoding phase will have 1 parameter being input string and will return a dictionary

```
''' Compress_LZ78  
    ==> input : input string  
    ==> output : dictionary & bit of this dictionary  
    ==> e.g: [[0, 'w'], [0, 'a'], [0, 'b'], [3, 'a'], [0, '_'], [1, 'a'], [3, 'b'], [2, '']]  
'''  
def Compress_LZ78(input):  
    compress_dict = []          ### init a dict with empty string  
    word = ''                  ### empty char  
    i = 0                      ### index  
    for key in input:          ### iterate input  
        i += 1  
        word += key            ### append a string  
        if not word in compress_dict: ## word is not in dict  
            compress_dict.append(word) ## append compress_dict  
            comp_dict.append([compress_dict.index(word[:-1]), word[-1]])  
            word = ''  
        elif i == len(input): ## if yes  
            comp_dict.append([compress_dict.index(word), ''])  
            word = ''  
        else:  
            pass  
    return comp_dict, bitMapping(comp_dict)
```

J'ai déjà commenté certaines lignes de code pour clarifier le flux de contrôle

I wrote some comments to clarify the workflow

LA STRUCTURE DE DONNÉE QUE J'AI UTILISÉ POUR SAUVEGARDER LE RÉSULTAT DE LA COMPRESSION EST DICTIONNAIRE

The data structure which I used to store the compression result is **DICTIONARY**

ALGORITHME DE LZ78 – PHASE DU DÉCODAGE

LZ78 – DECODING PHASE

Contrairement au phase d'encodage, la phase du décodage va reconstruire les informations originales. Pour cette phase, je vais utiliser 2 structures de données suivantes pour sauvegarder les résultats :

Unlike the encoding phase, the decoding phase will reconstruct the original information. For this phase, I will use the following 2 data structures to store the results:

1. Le dictionnaire APRES décodage Dictionary after decoding
2. La chaîne caractère originale Original string

```
''' Decompress_LZ78
    ==> input: input string as compressed string (obtained from Compress_LZ78)
    ==> output : 2 params ==> dictionary & original string
    ==> e.g : ... & wabba_wabba_wabba_woo_woo
'''
def Decompress_LZ78(compressedInput):
    p = ""
    outString = ""
    for (i, j) in compressedInput:      ### compressedInput is a dict <key, value>
        if i is not 0:                  ### if dict is exist -> recursive the decompression
            p = decomp_dict [i-1]
        else:                           ### empty dict
            p = ""
        decomp_dict.append(p + j)
    for s in decomp_dict :              ### concat string to have original string
        outString += s
    return dictMapping(bitMapping(decomp_dict)), outString      ### return a dict and original string
```

J'ai déjà commenté certaines lignes de code pour clarifier le flux de contrôle

I wrote some comments to clarify the workflow

LA STRUCTURE DE DONNÉE QUE J'AI UTILISÉ POUR SAUVEGARDER LE RÉSULTAT DE LA COMPRESSION EST DICTIONNAIRE et CHAÎNE DE CARACTÈRE ORIGINALE

The data structure which I used to store the compression result is **DICTIONARY** and **STRING**

En outre, je vais implémenter 2 méthodes pour convertir un dictionnaire en forme binaire (et vice-versa) en utilisant le package json

```
''' bitMapping ==> convert a dict to bit (0,1)'''  
def bitMapping(inputDict):  
    str = json.dumps(inputDict)  
    binary = ' '.join(format(ord(letter), 'b') for letter in str)  
    return binary  
  
''' dictMapping ==> convert bit to dict'''  
def dictMapping(inputBin):  
    jsn = ''.join(chr(int(x, 2)) for x in inputBin.split())  
    d = json.loads(jsn)  
    return d
```


TESTING L'ALGORITHME LZ78

LZ78 – TEST CASES

Les test-cases simples (le code Python ci-dessous) >>> SUPPRIMER LES COMMENTAIRES POUR LANCER LES TESTS

```
# #### simple test ==> please uncomment to test

# ## Uncomment to test
# data = 'wabba_ wabba_ wabba_ wabba_ wabba_ wabba_ wabba_ woo_woo_woo_woo'
# print("original_string_is : ", data)
# print("*****")

# # # # compress data ==> return a dict
# compress_data,bin_comp_data = Compress_LZ78(data)
# print("compressed_dict : ", compress_data)
# print("compressed_data_by_binary : ", bin_comp_data)

# # # # decompress data ==> return a dict and original string
# print("*****")
# decompress_data,original_string = Decompress_LZ78(compress_data)
# print("decompressed_dict : ", decompress_data)
# print("original_string : ", original_string)

# # # #### compress ratio = original_size / compress_size
# print("compress_ratio : ", (len(compress_data)/len(data) * 100), "%")

# with open("testcase1", "w") as f:
#     f.write("original_string_is : \r\n" + data)
#
# f.write("*****")
#
#     f.write("compressed_dict : \r\n" + str(compress_data))
#     f.write("compressed_data_by_binary : \r\n" + str(bin_comp_data))
#
# f.write("*****")
#
#     f.write("decompressed_dict : \r\n" + str(decompress_data))
#     f.write("original_string : \r\n" + str(original_string))
#
#     f.write("compress_ratio : " + str((len(compress_data)/len(data) * 100)) + "%")
```


Je sauvegarde les résultats dans les fichiers .txt (testcase1, testcase2, testcase3, testcase4)

1. wabba wabba wabba wabba woo woo woo (RATIO: 48,571....%)

```
original_string_is : wabba_wabba_wabba_wabba_woo_woo_woo
*****
compressed_dict : [[0, 'w'], [0, 'a'], [0, 'b'], [3, 'a'], [0, '_'], [1, 'a'], [3, 'b'], [2, '_'], [6, 'b'], [4, '_'], [9, 'b'], [8, 'w'], [0, 'o'], [13, '_'], [1, 'o'], [14, 'w'], [13, 'o']]
compressed_data_by_binary : 1011011 1011011 110000 101100 100000 100010 1110111 100010 1011101 101100 100000 1011011 110000 101100 100000 100010 100001 100010 1011101 101100 100000 1011011
110000 101100 100000 100010 1100010 100010 1011101 101100 100000 1011011 11001 101100 100000 100010 1100001 100010 1011101 101100 100000 1011011 110000 101100 100000 100010 101111 100010 1011101
101100 100000 1011011 110001 101100 100000 100010 100000 100010 1100001 100001 1011011 101100 100000 1011011 101100 100000 1011011 101100 100000 1011011 101100 100000 1011011
1011111 100010 1011101 101100 100000 1011011 10110 101100 100000 100010 101101 101100 100000 100010 101111 100010 1011101 101100 100000 1011011 101100 100000 100010 101111 100010 1011101
111001 101100 100000 100010 1100010 100010 1011101 101100 100000 1011011 111000 101100 100000 100010 1101111 100010 1011101 101100 100000 1011011 110000 101100 100000 100010 1101111 100010 1011101
101100 100000 1011011 110001 110011 101100 100000 100010 1011111 100010 1011101 101100 100000 1011011 110001 101100 100000 100010 101111 100010 1011101 101100 100000 1011011 110001 110100 101100
100000 100010 1101111 100010 1011101 101100 100000 1011011 110001 101100 100000 100010 101111 100010 1011101 101100 100000 1011011 101100 100000 1011011
*****
decompressed_dict : ['w', 'a', 'b', 'ba', '_', 'wa', 'bb', 'a_', 'wab', 'ba_', 'wabb', 'a_w', 'o', 'o_', 'wo', 'o_w', 'oo']
original_string : wabba_wabba_wabba_wabba_woo_woo_woo
compress ratio : 48.57142857142857 %
```

2. wabba wabba wabba wabba wabba woo woo woo woo (RATIO: 45,652...%)

```

original_string_is : wabba_wabba_wabba_wabba__wabba_woo_woo_woo_woo
*****
compressed_dict : [[0, 'w'], [0, 'a'], [0, 'b'], [3, 'a'], [0, '_'], [4, 'a'], [3, 'b'], [2, '_'], [6, 'b'], [4, '_'], [9, 'b'], [8, '_'], [11, 'a'], [5, 'w'], [0, 'o'], [15, '_'], [1, 'o'], [16, 'w'], [15, 'o'], [14, 'o'], [15, '']]
compressed_data_by_binary : 0110111 1011011 110000 101100 100000 100010 1101011 100010 1011101 101100 100000 1011011 110000 101100 100000 100010 1100001 100010 1011101 101100 100000 1011011
110000 101100 100000 100010 1100010 100010 1011101 101100 100000 1011011 110011 101100 100000 100010 1100001 100010 101101 101100 100000 1011011 110000 101100 100000 100010 101111 100010 1011101
101100 100000 1011011 110001 101100 101100 100000 100010 1100001 100010 101101 101100 100000 1011011 110010 100000 1011011 100010 101100 100000 100010 1011011 100010 101100 100000 100010
110111 100010 100010 101101 101100 100000 1011011 100110 101100 100000 100010 100010 100000 100010 101101 101100 100000 100010 101101 101100 100000 100010 101101 101100 100000 100010
111001 101100 100000 100000 100010 100010 101101 101100 100000 100010 111000 101100 100000 100010 101111 100010 101101 101100 100000 1011011 110001 110001 101100 100000 100010 1100001 100010
101101 101100 100000 1011011 110101 101101 101100 100000 1011011 110111 100010 1101101 100000 101100 100000 100010 101111 100000 101101 101100 100000 100010 101101 101100 100000 100010
100000 100010 1011111 100010 101101 101100 100000 1011011 110001 101100 100000 100010 110111 100010 101101 101100 100000 100010 110111 100010 101101 101100 100000 100010 110101 101100
100000 1011011 110001 101101 101100 100000 100010 110111 100010 101101 100000 100010 110111 100010 101101 101100 100000 100010 110111 100010 101101 101100 100000 100010 110101 101100
100000 100010 100010 101101 101101
*****
decompressed_dict : ['w', 'a', 'b', 'ba', '_', 'wa', 'bb', 'a_', 'wab', 'ba_', 'wabb', 'a_', 'wabba', 'w', 'o', 'o_', 'wo', 'o_w', 'oo', 'w_o', 'o']
original_string : wabba_wabba_wabba_wabba__wabba_woo_woo_woo_woo
compress ratio : 45.65217391304348 %

```

3. **wabba wabba wabba wabba wabba wabba woo woo woo woo woo (RATIO: 42,857...%)**

[illegible]

4. wabba wabba wabba wabba wabba wabba wabba woo woo woo woo (RATIO: 44,44...%)

```
original_string_is : wabba_wabba_wabba_wabba_wabba_wabba_woo_woo_woo_woo
*****
compressed_dict : [[0, 'w'], [0, 'a'], [0, 'b'], [3, 'a'], [0, '_'], [0, ''], [1, 'a'], [3, 'b'], [2, '_'], [6, 'w'], [2, 'b'], [4, '_'], [10, 'a'], [8, 'a'], [5, ''], [7, 'b'], [12, ''], [16, 'b'], [9, ''], [18, 'a'], [5, 'w'], [0, 'o'], [22, '_'], [2, ''], [1, 'o'], [23, 'w'], [22, 'o'], [21, 'o'], [22, '']]
compressed_data_by_binary : 1011011 1011011 100000 101100 100000 100000 1101111 000010 1011101 101100 100000 1011011 110000 101100 100000 100010 1100001 100010 1011101 101100 100000 1011011
101000 101100 100000 100010 1100001 100000 1011011 101100 100000 100010 1011011 101100 100000 100000 1011011 100000 101100 100000 100010 1011111 000010 1011011
101100 100000 1011011 110000 101100 100000 100000 100000 1011011 101100 100000 1011011 110001 101100 100000 100010 1011011 101100 100000 1011011 110011 101100 100000 100010
1100001 100010 1011011 101100 100000 1011011 101100 100000 100010 1011111 000010 1011011 101100 100000 100010 1011011 101100 100000 100010 1101011 101100 100000 100010
1010001 100010 1011011 101100 100000 1011011 101100 100000 100010 1011111 000010 1011011 101100 100000 100010 1011011 101100 100000 100010 1101011 101100 100000 100010
1010001 100010 1011011 101100 100000 100010 100000 100010 1011011 101100 100000 1011011 110110 101100 100000 100010 1101011 101100 100000 100010 1101011 101100 100000 100010
100000 1011011 110001 101100 100000 100010 100000 100010 1011011 101100 100000 100010 110001 110000 101100 100000 100000 1011011 101100 100000 100010 1011011 101100 100000 100010
1011011 100000 1011011 101100 100000 1011011 100000 101100 100000 100010 1011011 000010 1011011 101100 100000 100010 1011011 101100 100000 100010 1011011 101100 100000 100010
1011011 100001 101100 100000 100010 1011011 100000 1011011 101100 100000 100010 1011011 100000 100000 100010 1011011 100000 100000 100010 1011011 100000 100000 100010
1010111 000010 1011011 101100 100000 100010 1011011 100000 100000 100010 1011011 000010 1011011 101100 100000 100010 1011011 101100 100000 100010 1011011 101100 100000 100010
1010111 000010 1011011 101100 100000 100010 1011011 100000 100000 100010 1011011 000010 1011011 101100 100000 100010 1011011 101100 100000 100010 1011011 101100 100000 100010
1010111 000010 1011011 101100 100000 100010 1011011 100000 100000 100010 1011011 000010 1011011 101100 100000 100010 1011011 101100 100000 100010 1011011 101100 100000 100010
1010111 000010 1011011 101100 100000 100010 1011011 100000 100000 100010 1011011 000010 1011011 101100 100000 100010 1011011 101100 100000 100010 1011011 101100 100000 100010
*****
decompressed_dict : ['w', 'a', 'b', 'ba', '_', ' ', 'wa', 'bb', 'a', ' ', 'w', 'ab', 'ba_', 'wa', 'bba', ' ', 'wab', 'ba_', 'wabb', 'a ', 'wabba', 'w', 'o', 'o_', 'wo', 'o_w', 'oo', 'w_o', 'o']
original_string : wabba_wabba_wabba_wabba_wabba_wabba_woo_woo_woo_woo
compress_ratio : 44.44444444444444 %
```


Les test-cases complexes (fichiers texte) – Le code Python ci-dessous

```
# ### complex test with text file
file = "text.txt"
compress_data = []

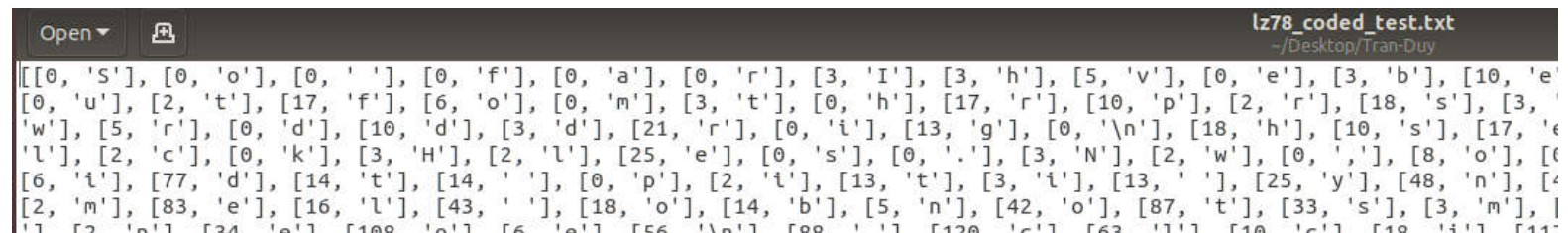
with open(file, "r") as file:
    data = file.read()
    data_size = len(data)
    _, bin_comp_data = Compress_LZ78(data)

    with open("lz78_coded_test.txt", "w") as f:
        compress_data, _ = Compress_LZ78(data)
        f.write(str(compress_data))

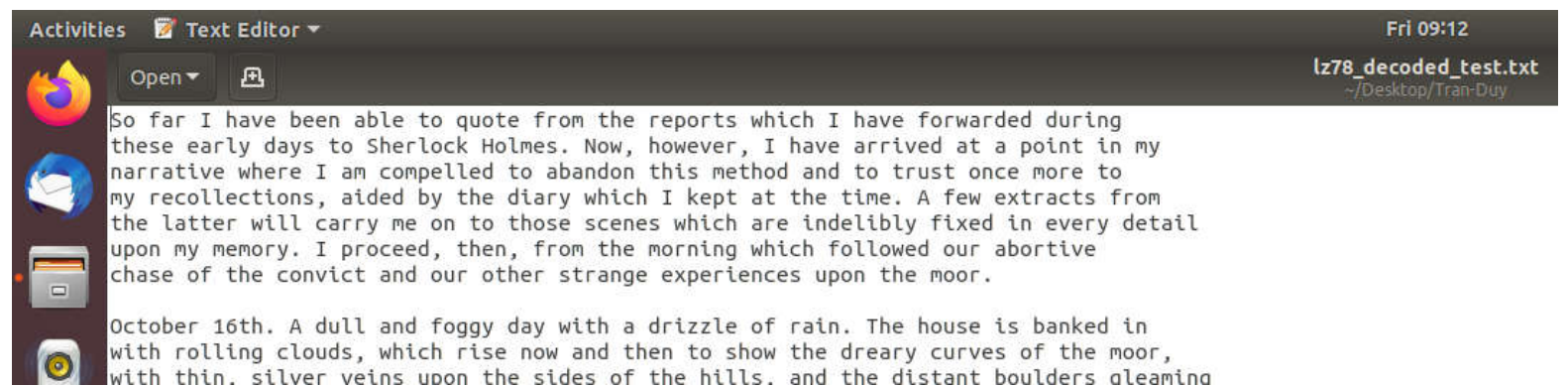
    with open("lz78_decoded_test.txt", "w") as f:
        _, decompress_data = Decompress_LZ78(compress_data)
        f.write(str(decompress_data))
```

J'ai déjà sauvegarder les résultats dans 3 fichiers texte :

1. « lz78_code.txt » sauvegarde le dictionnaire d'encodage
2. « lz78_decoded.txt » sauvegarde les informations originaux
3. « testcase5.txt » sauvegarde la ratio de compression



lz78_coded.txt



lz78_decoded.txt



testcase5.txt

ratio is : 50.96472320067294%

La ratio de compression

J'ai fait une autre test pour un autre texte (chapitre 1 – La Dame au Camelia) [je l'ai mis ce fichier dans la même répertoire). Les résultats que je vais obtenir sont :

Activities

Text Editor

Fri 10:50

lz78_coded_dame.txt

~/Desktop/Tran-Duy

[[0, 'I'], [0, 'n'], [0, ' '], [0, 'm'], [0, 'y'], [3, 'o'], [0, 'p'], [0, 'i'], [2, 'i'], [0, 'o'], [2, ' '], [3, 'e'], [3, 't'], [10, ' '], [0, 'c'], [0, 'r'], [21, 'a'], [13, 'e'], [3, 'c'], [0, 'h'], [0, 'a'], [25, 'a'], [34, 's'], [7, 'e'], [36, '\n'], [30, ' '], [20, 'o'], [2, 'g'], [22, 'i'], [4, 'e'], [12, 'n'], [3, 's'], [8, 's'], [15, 'p'], [10, 's'], [18, 'b'], [20, 'e'], [22, 'o'], [50, 'p'], [26, 'k'], [57, ' '], [20, 'a'], [3, 'b'], [21, 'e'], [2, ' '], [17, 'e'], [25, 'i'], [10, 'u'], [17, 'l'], [5, ' '], [30, 'c'], [0, 'q'], [1, 'l'], [52, ' '], [21, 'n'], [84, 'g'], [29, ' '], [13, 'o'], [49, 'v'], [99, 't'], [56, ' '], [1, '\n'], [24, 'o'], [31, 't'], [95, 'g'], [105, 'a'], [2, 'd'], [3, 'I'], [94, 'g'], [22, 'h'], [39, 'r'], [26, 'd'], [33, ' '], [13, 'h'], [39, 't'], [25, 'u'], [135, ' '], [10, 'f'], [68, 's'], [102, 'r'], [86, 'i'], [81, 'w'], [29, 'i'], [105, 'w'], [114, ' '], [135, 'e'], [3, 'e'], [0, 'x'], [24, 'e'], [7, 't'], [8, 'o'], [81, 'o'], [112, '\n'], [1

Activities

Text Editor

Fri 10:51

lz78_decoded_dame.txt

~/Desktop/Tran-Duy

In my opinion, it is impossible to create characters until one has spent a long time in studying men, as it is impossible to speak a language until it has been seriously acquired. Not being old enough to invent, I content myself with narrating, and I beg the reader to assure himself of the truth of a story in which all the characters, with the exception of the heroine, are still alive. Eye-witnesses of the greater part of the facts which I have collected are to be found in Paris, and I might call upon them to confirm me if my testimony is not enough. And, thanks to a particular circumstance, I alone can write these things, for I alone am able to give the final details, without which it would have been impossible to make the story at once interesting and complete.

Activities

Text Editor

Fri 10:51

testcase6.txt

~/Desktop/Tran-Duy

ratio is : 56.428422588039176%

CONCLUSION

L'algorithme de LZ78 que j'ai déjà implémenté a marché bien en 2 cases : simple et complexe. **Cependant, je pense que ces algorithmes ne sont pas très efficace (la complexité est grande car j'ai déjà utilisé la récursive) => la difficulté est comment annuler la récursive (je pense pour faire cela, je peux utiliser l'approche de l'algorithme de gourmand).**

UNE PETITE OBSERVATION POUR LES TEST CASES

1. POUR LES CAS DE TEST SIMPLE : LA RATIO EST ENTRE 40 – 50%
2. POUR LES CAS DE TEST COMPLEXE : LA RATIO EST ENTRE 50 – 60%