

Belhadj Walid

Master II Sicom

Travail pratique : Sécurité des cartes à puces

Sujet : codage Huffman (compression et décompression)

Explication et démarche :

Code huffman en un fichier **main.py**

Exécution : python main.py

Cette exécution retourne la compression et décompression du fichier **textfile.txt**

Résultat de sortie : un fichier binaire (**textfile.bin**) et un fichier **textfile_decompressed.txt** qui doit être comparé avec le fichier original le code est suffisamment commenté,

Blocks du code :

Bibliothèques :

```
import json # stocke le chemin vers le fichier source
from heapq import heappop, heappush # gestion de files
from bitstring import BitArray #conversion binaire/bytes/string
import os #pour s'adapter au OS
```

Ouverture des fichiers nécessaire :

```
# on ouvre et on charge le fichier à partir du fichier config.json
config = open('./config.json')
config_json = json.load(config)
# renvoie une liste contenant chaque ligne du fichier en tant qu'élément
# de liste
# utilise le paramètre d'encodage pour ouvrir le fichier et le lire
# dans n'importe quelle langue possible
file = open(config_json["filepath_text"], encoding="utf8").readlines()
# print(file)
```

Fonction de gestion de fréquence d'apparition des caractères :

```
def count_frequencies(file) -> dict:
    """
    Cette fonction ouvre un fichier texte en argument puis renvoie le
    fréquences chaque caractère apparaît dans le fichier texte. Les personnages
    sont ensuite triés par ordre en partant de la fréquence la plus basse jusqu'à
    le plus haut. Valeurs triées par ordre croissant.
    """
```

La sortie :

```
# stocker le dictionnaire des fréquences dans une variable
# appelée "d" pour une utilisation ultérieure
d = count_frequencies(file)
print(d)
```

Définition de l'arbre de Huffman :

```
def creation_of_huffmantree(d) -> list:
    """
    Cette fonction obtient le dictionnaire renvoyé par la fonction
    count_frequencies
    et renvoie une liste de listes qui représentent l'arbre de
    Huffman des personnages
    dans le fichier texte. Il est formaté de la manière suivante :
    [[character_frequency, [character, code]]
    Au début, le code sera une chaîne vide mais il aura du contenu plus tard.
    # affecte une variable aux listes de listes pour la représentation
    # de l'arbre de Huffman
    # utilise la méthode items() pour retourner un objet de vue
    # L'objet view contient les paires clé-valeur du dictionnaire,
    # sous forme de tuples dans une liste.
    huffman_tree = [[frequency, [char, ""]] for char, frequency in d.items()]
    # on utilise la boucle while pour créer un arbre huffman
    # tant que la longueur de la liste de l'arbre de Huffman est supérieure à 1
    while len(huffman_tree) > 1: ...
    return huffman_list
print(creation_of_huffmantree(d))
```

Fonction d'encodage du textes (entrée/sortie)

```
def encoded_texts(huffman_list: list, new_list: list) -> str:
    """
    Cette fonction obtient la liste renvoyée par la fonction
    creation_of_huffmantree
    et une nouvelle liste comme paramètre aussi. Cette fonction permet
    d'obtenir le code
    texte, il remplace les caractères par le contexte du fichier
    texte d'origine par
    codes attendus et renvoie donc un texte encodé.
    """
```

```
    encoded_text = string.translate(convert)
    # print(encoded_text)
    return dictionary, encoded_text
```

Fonction de remplissage de textes :

```
120
121 def padding_text(encoded_text: str):
122
123     """
124     Cette fonction obtient la chaîne encoded_text de la fonction
125     encoded_text ci-dessus et renvoie la version complétée du texte encodé.
126     Ajoute des caractères au format texte codé tel qu'il doit être affiché.
127     Ajoute des bits à la chaîne encodée donc la longueur
128     est un multiple de 8 et peut donc être encodé efficacement.
129     """
```

Return :

Elle joint le texte rembourré et le texte encodé pour obtenir la version finale du texte encodé **encoded_text = padded_data + encoded_text**

return encoded_text

```
def compression():
    """
    Cette fonction est formatée pour être utilisée lors de la compression
    le fichier texte. Il utilise les valeurs renvoyées dans les
    fonctions pour obtenir un fichier correctement compressé :
    --> fonction count_frequencies
    --> fonction creation_of_huffmantree
    --> fonction encoded_text
    --> fonction padding_text
    Tout ce qui précède
    """
```

L'écriture se fait sur un fichier binaire (.bin)

```
return output_path
```

Fonction de suppression de bourrage :

```
4
5 def remove_padding(bit_string):
6
7     """
8     Cette fonction est utilisée dans la fonction de décompression pour
9     décompresser le fichier texte. Sa fonctionnalité est de supprimer
10    le remplissage ajouté à la fonction de remplissage de texte pour
11    obtenir une décompression correcte du texte.
12    Elle renvoie le texte encodé sans le remplissage.
13    """
14
15    padded_data = bit_string[:8]
16    extra_padding = int(padded_data, 2)
17
18    bit_string = bit_string[8:]
19    encoded_text = bit_string[:-1*extra_padding]
20
21    return encoded_text
```

Fonction de décompression :

```
def decompress(input_path):
    """
    Cette fonctionnalité de fonctions est de décompresser le déjà compressé
    fichier texte d'avant. Il utilise les valeurs renvoyées dans les
    fonctions pour obtenir un fichier correctement compressé :
    --> fonction creation_of_huffmantree
    --> fonction encoded_text
    --> fonction remove_padding
    Tout ce qui précède
    """
```

Fonction de décodage de bytes :

```
def decode_text(encoded_text: str, reverse_mapping: dict):
    """
    Cette fonction renvoie le texte décodé du fichier texte.
    Utilise les boucles for et l'instruction if pour obtenir le texte décodé.
    """
```

Partie opérations :

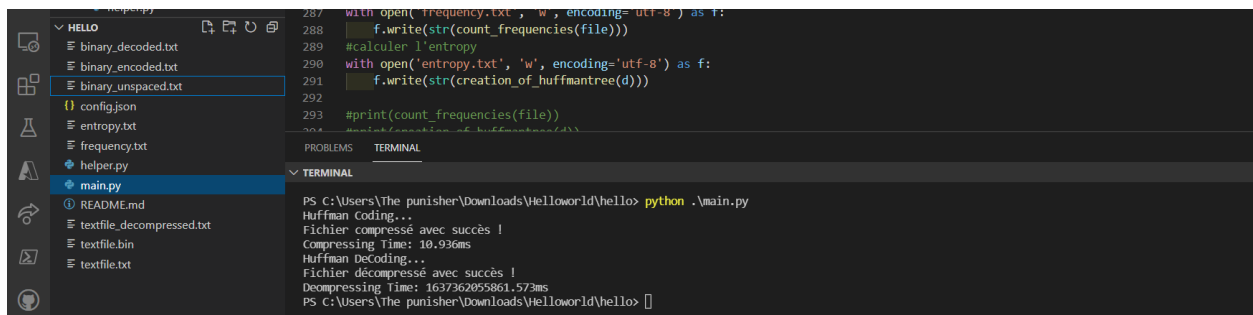
```
print("Huffman Coding...")
timeofcompress = time.time()
# Compter l'apparition de chaque caractère
with open('frequency.txt', 'w', encoding='utf-8') as f:
    f.write(str(count_frequencies(file)))
#calculer l'entropy
with open('entropy.txt', 'w', encoding='utf-8') as f:
    f.write(str(creation_of_huffmantree(d)))

#print(count_frequencies(file))
#print(creation_of_huffmantree(d))
#stock binary text
with open('binary_encoded.txt', 'w', encoding='utf-8') as f:
    #print(encoded_texts(file, creation_of_huffmantree(d))[1])

with open('binary_decoded.txt', 'w', encoding='utf-8') as f:
    #print(padding_text(encoded_texts(file, creation_of_huffmantree(d))[1]))

with open('binary_unspaced.txt', 'w', encoding='utf-8') as f:
    print("Fichier compressé avec succès ! ")
timeofcompress = time.time()-timeofcompress
print('Compressing Time: ', end='')
print(str(round(timeofcompress * 1000, 3)) + 'ms')
#print(compression())
print("Huffman DeCoding...")
timeofdecompress = time.time()
decompress(compression())
#print(decompress(compression()))
print('Decompressing Time: ', end='')
print(str(round(timeofdecompress * 1000, 3)) + 'ms')
```

Résultat et fichiers créés :



Amélioration :

Comparaison des deux fichiers :

```

328 #comparaison
329 print("Phase de comparaison")
330 > with open('textfile.txt', 'r') as file1: ...
333
334 check=False
335 same.discard('\n')
336 if same != check:
337 |
338 |     print("les fichiers sont identiques")
339 | else:
340 |     print("les fichiers ne sont pas identiques")
341
342 with open('save_difference.txt', 'w') as file_out:
343 |     for line in same:
344 |         file_out.write(line)
345

```

Enregistrements de la table de fréquences, binary files ...

```

286 print("Huffman Coding...")
287 timeofcompress = time.time()
288 # Compter l'apparition de chaque caractère
289 with open('frequency.txt', 'w', encoding='utf-8') as f:
290 |     f.write(str(count_frequencies(file)))
291 #calculer l'entropy
292 > with open('entropy.txt', 'w', encoding='utf-8') as f: ...
297 #stock binary text
298 > with open('binary_encoded.txt', 'w', encoding='utf-8') as f: ...
300 #print(encoded_texts(file, creation_of_huffmantree(d))[1])
301 > with open('binary_decoded.txt', 'w', encoding='utf-8') as f: ...
303 #print(padding_text(encoded_texts(file, creation_of_huffmantree(d))[1]))
304 > with open('binary_unspaced.txt', 'w', encoding='utf-8') as f: ...
305

```

Compter le temps de compression et décompression

```

310 print('Compressing Time: ', end='')
311 print(str(round(timeofcompress * 1000, 3)) + 'ms')
312 inputsize=os.path.getsize("textfile.txt")
313 outputsize=os.path.getsize("textfile.bin")
314 compression_ratio=float(outputsize/inputsizes)*100
315

```

```

318 print("\n")
319 print("Huffman DeCoding...")
320 timeofdecompress = time.time()
321 decompress(compression())
322 #print(decompress(compression()))
323 print("\n")
324 print('Deompressing Time: ', end='')
325 print(str(round(timeofdecompress * 1000, 3)) + 'ms')
326

```

Résultat final avec le fichier original

HELLO

binary_decoded.txt

binary_encoded.txt

binary_unspaced.txt

config.json

entropy.txt

frequency.txt

main.py

save_difference.txt

textfile_decompressed.txt

textfile.bin

textfile.txt

white_paper.pdf

PS C:\Users\The punisher\Downloads\Helloworld\hello> python .\main.py

Huffman Coding...

Fichier compressé avec succès !

Compressing Time: 4072.449ms

Compression percentage (Compressed size/file size) (%): 55.55630377259836

Huffman DeCoding...

Fichier décompressé avec succès !

Decompressing Time: 1637371588052.772ms

Phase de comparaison

les fichiers sont identiques

PS C:\Users\The punisher\Downloads\Helloworld\hello>