

CORRECTION DU TD itc² N° 6: DICTIONNAIRES 1: UTILISATIONS ÉLÉMENTAIRES, TABLE DE HACHAGE, COLLISIONS

EXERCICE N°1:

Algorithme de compression par dictionnaire: l'algorithme

LZ78

- ❶ On propose le code complété ainsi:

Listing 1: Fonction compresser(texte,code,dc,i)

```

1 def compresser(texte ,code ,dc ,i):
2     n=len(texte)
3     assert i<len(texte)
4     w=""
5     while i<n and w+texte[i] in dc: #on avance sur le préfixe
6         w+=texte[i]
7         i+=1
8     if i<n:
9         ###le préfixe w est dans le dictionnaire###
10        ###le nouveau préfixe wl=w+texte[i] n'y figure pas et y
11        est ajouté###
12        p=len(dc) #la position est donnée par la longueur du
13        dictionnaire
14        c=texte[i]
15        wl=w+c
16        dc[wl]=p
17        code.append((dc[w],c))
18    return code , i+1

```

- ❷ On obtient à partir du texte "ABRACADABRA":

$$code = [(0, 'A'), (0, 'B'), (0, 'R'), (1, 'C'), (1, 'D'), (1, 'B'), (3, 'A')]$$

$$dc = \{": 0, 'A' : 1, 'B' : 2, 'R' : 3, 'AC' : 4, 'AD' : 5, 'AB' : 6, 'RA' : 7\}$$

- ❸ Le code pour inverser clés et valeurs est élémentaire en procédant par compréhension:

$dc_i = \{nvcl: nvval \text{ for } (nvval, nvcl) \text{ in } dc.items()\}$

Le dictionnaire inversé donne pour le texte "ABRACADABRA":

$$dc_i = \{0 : "", 1 : 'A', 2 : 'B', 3 : 'R', 4 : 'AC', 5 : 'AD', 6 : 'AB', 7 : 'RA'\}$$

- ❹ Principe de reconstruction du texte:

On rappelle les contenus de code, dc et dc_i :

$$code = [(0, 'A'), (0, 'B'), (0, 'R'), (1, 'C'), (1, 'D'), (1, 'B'), (3, 'A')]$$

$$dc = \{": 0, 'A' : 1, 'B' : 2, 'R' : 3, 'AC' : 4, 'AD' : 5, 'AB' : 6, 'RA' : 7\}$$

$$dc_i = \{0 : "", 1 : 'A', 2 : 'B', 3 : 'R', 4 : 'AC', 5 : 'AD', 6 : 'AB', 7 : 'RA'\}$$

- On inverse le dictionnaire dc i.e. les valeurs (str) et les clés (int).
- On initialise un texte vide `texte=""`.
- On itère sur code (variable `e1`);
 - le premier élément du tuple `e1[0]` i.e. la position du caractère est extrait et permet de récupérer dans le dictionnaire inversé le préfixe avant le caractère qui va être traité: `dc_i[e1[0]]`.
 - On récupère ensuite le caractère à ajouter après le préfixe avec `e1[1]`.
 - On concatène enfin `texte`, le préfixe avant le caractère à ajouter `dc_i[e1[0]]` et le caractère traité `e1[1]`
- ❺ Le code correspondant est le suivant:

Listing 2:

```
1 def decompression(code , dc):
2     dci={ nvcle: nvval for (nvval , nvcle) in dc.items() } #
3     construction du dictionnaire inversé
4     texte=""
5     print( " Dictionnaire inversé:" , dci)
6     for el in code:
7         pref=dci[el[0]]
8         texte+=pref+el[1]
9     return texte
```

Le déroulement pas à pas du code donne:

el	w	s	texte
(0,"A")	""	"A"	"A"
(0,"B")	""	"B"	"AB"
(0,"R")	""	"R"	"ABR"
(1,"C")	"A"	"C"	"ABRAC"
(1,"D")	"A"	"D"	"ABRACAD"
(1,"B")	"A"	"B"	"ABRACADAB"
(3,"A")	"R"	"A"	"ABRACADABRA"

EXERCICE N°2:

Fonction de hachage - résolution d’une collision

❶ On obtient les hachages suivants:

k	h(k)
26	→ 0
37	→ 11
24	→ 11
30	→ 4
11	→ 11

❷ La table de hachage comporte 13 alvéoles numérotées de 0 à 12. Après résolution des collisions par adressage ouvert et sondage linéaire, on obtient l’occupation suivante de la table:

$h_i(k) = (h(k) + i) \% m$	k
0	26
1	11 (3 collisions)
2	
3	
4	30
5	
6	
7	
8	
9	
10	
11	37
12	24 (1 collision)

❸ On constate que la taille de la table va être insuffisante pour ajouter 1,2,3,4,5,6,7,8,9,10 :

$h_i(k) = (h(k) + i) \% m$	k
0	26
1	11 (3 collisions)
2	1 (1 collision)
3	2 (1 collision)
4	30
5	3 (2 collisions)
6	4 (2 collisions)
7	5 (2 collisions)
8	6 (2 collisions)
9	7 (2 collisions)
10	8 (2 collisions)
11	37
12	24 (1 collision)

Les valeurs 9 et 10 ne trouvent plus d’alvéoles libres avec l’adressage ouvert et sondage linéaire.

Deux solutions sont envisageables:

- On augmente la taille du dictionnaire (c’est à dire on augmente m).
- On procède à un chainage (vu en cours).

EXERCICE N°3:

Exploration détaillée d’une collision

Listing 3:

```
1. 1 def h(k):
2     res=ord(k[0])*256+ord(k[1]) #on initialise la valeur de res
   avant la boucle
3     for i in range(2, len(k)):
4         res=res*256+ord(k[i])
5     return res
```

Autre proposition, toujours itérative:

Listing 4:

```
1 def h_2(ch):
2     for i in range(0, len(ch)-1):
3         if i==0:
4             res=ord(ch[i])*256+ord(ch[i+1])
5         else:
6             res=res*256+ord(ch[i+1])
7     return res
```

Enfin, le schéma de Horner est particulièrement adapté à une approche récursive:

- la construction de la clé hachée se fait cette fois en partant de la fin (coefficient a_0 calculé en premier); le code renverra le cas récursif avec:

`return ord(ch[-1])+256*h_rec(ch[:-1])`

- le cas de base va intervenir lorsqu'il restera uniquement dans la chaîne les deux premiers caractères, correspondant aux deux coefficients d'ordre plus élevé a_n et a_{n-1} :

`return ord(ch[0])*256+ord(ch[1])`

Cela donne:

Listing 5:

```
1 def h_rec(ch): #on part cette fois du dernier caractère
2     if len(ch)==2: #cas de base
3         return ord(ch[0])*256+ord(ch[1])
4     else: #cas récursif
5         return ord(ch[-1])+256*h_rec(ch[:-1])
```

Le code précédent ne fonctionne qu'avec des clés de longueur minimale 2; pour corriger cela et envisager des clés à un caractère on peut faire:

Listing 6:

```
1 def h_rec(ch): #on part cette fois du dernier caractère
2     if len(ch)==1: #cas de base
3         return ord(ch)
4     else: #cas récursif
5         return ord(ch[-1])+256*h_rec(ch[:-1])
```

Listing 7:

```
2. 1 def entier_chaine(e):
   res=""
3     while e%256!=0:
4         res=chr(e%256)+res
5         e=e//256
6     return res
```

Du fait de la définition de la fonction h , chaque chaîne de caractère possède une image unique dans \mathbb{N} , la fonction h est donc injective.

Listing 8:

```
3. 1 def f(k):
   return h(ch)%255
```

4. On obtient la sortie suivante:

```
>>> print h("pouet")
47
>>> print h("chariot")
236
>>> print h("haricot")
236
```

5. La valeur obtenue par application du modulo sur l'entier "long" e peut s'écrire:

$$(a_n \cdot 256^n + a_{n-1} \cdot 256^{n-1} + \dots + a_1 \cdot 256 + a_0 \cdot 256^0)[255]$$

$$= (a_n + a_{n-1} + \dots + a_1 + a_0)[255]$$

L'indice renvoyé est donc indépendant de l'ordre des coefficients $(a_n, a_{n-1}, a_{n-2}, \dots, a_0)$, et sera donc **le même pour des mots anagrammes**.

6. Script de vérification de validité du dictionnaire pour $f(ch)$:

NB: cette procédure ne fonctionne que dans le cas d'un dictionnaire dont les clés sont des chaînes.

Listing 9:

```
1 def dico_valide(dict):
2     res=True
3     listecles=[cle for cle in dict.keys()]
4     for i in range(len(listecles)):
5         for j in range(i+1,len(listecles)):
6             if f(listecles[i])==f(listecles[j]):
7                 res=False
8     return res
```

Enfin, on propose cette variante qui évite un balayage total du dictionnaire par l'emploi de boucles conditionnelles:

Listing 10:

```
1 def dico_valide_bis(dict):
2     res=True
3     listecles=[cle for cle in dict.keys()]
4     i=0
5     while res and i<len(listecles)-1:
6         j=i+1
7         while res and j<len(listecles)-1:
8             res=f(listecles[i])!=f(listecles[j])
9             j=j+1
10        i+=1
11    return res
```

Cette dernière variante exploite une liste de 256 entiers (d'indice $\in [0, 255]$), tous initialisés à 0. Cette liste va recueillir, à chaque valeur d'indice, le nombre d'occurrences du hachage de clé qui a justement donné la valeur de cet indice (liste de fréquences); pour assurer l'absence d'anagrammes, aucun total d'occurrences ne doit dépasser 1: on teste donc cela directement dans la boucle conditionnelle de balayage de la liste des clés; la complexité est ainsi nettement réduite par rapport aux scripts proposés plus haut:

Listing 11:

```
1 def dico_valide_ter(dict):
2     res=True
3     listecles=[cle for cle in dict.keys()]
4     tabfreq=[0]*256
5     i=0
6     while res and i<len(listecles):
7         if tabfreq[f(listecles[i])]==0:
8             tabfreq[f(listecles[i])]+=1
9             i+=1
10        else:
11            res=False
12    return res
```