

TP n° 12 – Dictionnaires

On donne en Annexe quelques fonctions et méthodes liées aux dictionnaires python.

I Décompte des occurrences de lettres dans un mot ou un texte

D'après une idée de Serge Bays.

Exercice 1. On considère un mot, représenté en python par une chaîne de caractères (objet python de type `str`). Consulter l'Annexe du TP 03 pour un rappel sur les fonctions et méthodes sur les chaînes de caractères.

Dans un script python, écrire une fonction `comptage(mot)` qui prend comme argument une chaîne de caractères et renvoie un dictionnaire dont chaque clé est un caractère différent du mot et l'élément associé à la clé est le nombre de fois que ce caractère apparaît dans le mot.

Solution 1.

```

1 def comptage(mot):
2     '''renvoie un dictionnaire contenant le nombre
3     d'occurrences de chaque caractère du mot'''
4     occurrences = {}
5     for caractere in mot:
6         if caractere in occurrences:
7             occurrences[caractere] = occurrences[caractere] + 1
8         else:
9             occurrences[caractere]=1
10    return occurrences

```

Exercice 2. Tester la fonction précédente sur le mot de votre choix et en afficher le dictionnaire des occurrences des caractères. Puis la tester sur un deuxième mot contenant exactement les mêmes caractères que le premier mot, avec le même nombre d'occurrences pour chaque caractère, mais avec les caractères dans un ordre différent.

Tester si les dictionnaires des occurrences des caractères des deux des mots sont identiques. Conclusion ?

Solution 2.

```

1 print(comptage('tatayoyo'))
2 print(comptage('yoayttoa'))
3 print(comptage('tatayoyo') == comptage('yyaottoa'))

```

À l'affichage d'un dictionnaire, python affiche les couples clé/élément dans l'ordre dans lequel ils ont été ajoutés au dictionnaire. Mais pour python deux dictionnaires sont identiques s'ils sont constitués exactement des mêmes couples clé/élément, quel que soit l'ordre dans lequel on les a ajoutés. C'est une différence avec les listes, qui ne sont égales que si tous les éléments sont dans le même ordre.

Exercice 3. On considère maintenant un texte, qui est aussi représenté en python par une chaîne de caractères. La différence entre un mot et un texte provient du fait que dans un texte il y a des caractères supplémentaires : la ponctuation (point, virgule, point-virgule, deux points, point d'exclamation, point d'interrogation), des espaces et des caractères spéciaux dits d'échappement dont le plus courant est celui qui signifie le passage à la ligne (`\n`).

Définir une fonction `statistiques(texte)` qui prend comme argument une chaîne de caractères et renvoie un dictionnaire dont chaque clé est un caractère différent du texte, à l'exclusion des ponctuations, espaces et caractère d'échappement, et l'élément associé à la clé est le nombre de fois que ce caractère apparaît dans le texte.

Solution 3.

```

1 def statistiques(texte):
2     '''renvoie un dictionnaire contenant le nombre

```

```

3     d'occurrences de chaque caractère du texte
4     à l'exclusion des ponctuations, espaces et saut de ligne'''
5     occurrences = {}
6     for caractere in texte:
7         if caractere not in " ,;.:!?\\n":
8             if caractere in occurrences:
9                 occurrences[caractere] = occurrences[caractere] + 1
10            else:
11                occurrences[caractere]=1
12    return occurrences

```

Exercice 4. Copier dans votre répertoire personnel les fichiers hugo.txt et perec.txt présents dans le répertoire Ressources. Avec python : écrire les instructions pour ouvrir hugo.txt, récupérer le contenu du fichier sous la forme d'une chaîne de caractères unique, fermer le fichier. Appliquer la fonction `statistiques` au contenu de ce fichier et afficher à l'écran les statistiques des occurrences des caractères dans le texte.

Solution 4.

```

1 file=open('hugo.txt','r')
2 hugo=file.read()
3 file.close()
4 print(statistiques(hugo))

```

Exercice 5. On souhaite trier les données pour lister les caractères dans l'ordre décroissant de leurs occurrences dans un texte. La méthode `.items()` permet d'accéder à l'ensemble des couples clé/élément sous la forme d'un objet python dédié de type `dict_items`. Afficher l'ensemble de ces couples pour le contenu du fichier hugo.txt.

Dans le but de trier les données, il peut être intéressant de les manipuler sous forme de liste, ce qui permet de réviser les algorithmes de tri vus au TP 07. Pour transformer l'objet précédent en liste, il suffit d'utiliser la fonction `list`. Exemple :

```

1 dico = {'a' : 5, 'b' : 12} # un dictionnaire
2 couples = dico.items() # objet dict_items contenant les couples clé/élément
3 liste = list(dico.items()) # une liste des couples clé/élément

```

Écrire une fonction `tri_occurrences` qui prend en argument un dictionnaire des nombres d'occurrences des caractères d'un texte et renvoie une liste des couples clé/élément triés dans l'ordre décroissant de leurs occurrences dans un texte.

L'afficher pour le contenu du fichier hugo.txt puis pour celui du fichier perec.txt.

Quelle est la particularité du second texte ?

Solution 5.

```

1 print(statistiques(hugo).items())
2
3 #Tri par insertion
4 def tri_occurrences(dictionnaire):
5     '''Tri par insertion des données du dictionnaire
6     après obtention de la liste des couples clé/élément'''
7     # création de la liste des couples clé/élément
8     liste=list(dictionnaire.items())
9     # Parcours de 1 à la taille de la liste
10    for i in range(1, len(liste)):
11        # On mémorise la position initiale de l'élément
12        k = liste[i]
13        # On décale cet élément vers la gauche tant que
14        # l'élément précédent est plus grand que lui
15        j = i-1

```

```
16     while j >= 0 and k[1] > liste[j][1] :
17         liste[j + 1] = liste[j]
18         j -= 1
19     liste[j + 1] = k
20     return liste
21
22 def tri_occurrences2(dictionnaire):
23     '''Tri à bulles des données du dictionnaire
24     après obtention de la liste des couples clé/élément'''
25     # création de la liste des couples clé/élément
26     liste=list(dictionnaire.items())
27     # Parcours de 1 à la taille de la liste
28     for i in range(1, len(liste)):
29         # Parcours des éléments précédents
30         for j in range(0, len(liste)-i):
31             # On permutte les deux éléments successifs
32             if liste[j][1] < liste[j+1][1] :
33                 liste[j], liste[j+1] = liste[j+1], liste[j]
34     return(liste)
35
36 print(tri_occurrences(statistiques(hugo)))
37 file=open('perec.txt','r')
38 perec=file.read()
39 file.close()
40 print(tri_occurrences(statistiques(perec)))
```

Le texte contenu dans perec.txt ne contient aucun « e » qui est pourtant la lettre la plus utilisée dans la langue française. Il s'agit d'un extrait de *La Disparition* de Georges Pérec, roman intégralement écrit sans la lettre e.

II Anagrammes

D'après une idée de Mickaël Péchaud.

Une anagramme d'un mot est obtenue par permutation des lettres de ce mot. Dans une anagramme, on ignore la casse (majuscule/minuscule), les symboles diacritiques (accents et cédilles), tirets et autres apostrophes.

Ainsi, les mots suivants sont des anagrammes :

- « nectar » et « carnet » ;
- « sortie » et « rôties » ;
- « niche » et « Chine ».

Exercice 6. Pour éliminer tous les symboles diacritiques, tirets, apostrophes et passer toutes les lettres en minuscules, on fournit le code d'une fonction :

```

1 import unicodedata
2 def enleve_diacritiques_majuscules_tirets(mot):
3     '''renvoie le mot en minuscules, sans les diacritiques,
4     tirets, apostrophes, etc.'''
5     return unicodedata.unidecode(mot).lower().replace('-', '')

```

Écrire une fonction `sont_anagrammes` qui prend en argument deux mots, teste si les deux mots sont des anagrammes l'un de l'autre et renvoie le résultat du test, c'est-à-dire un booléen (True ou False).

Solution 6.

```

1 import unicodedata
2 def enleve_diacritiques_majuscules_tirets(mot):
3     '''renvoie le mot s en minuscule, et sans les accents et cedilles'''
4     return unicodedata.unidecode(mot).lower().replace('-', '')
5
6 def sont_anagrammes(mot1, mot2):
7     '''teste si deux mots sont anagrammes l'un de l'autre après
8     avoir éliminé les diacritiques, accents, etc.'''
9     return comptage(enleve_diacritiques_majuscules_tirets(mot1)) == comptage(enleve_diacritiques_ma
10
11 print(sont_anagrammes('tatayoyo', 'yôyôtâtâ'))
12 print(sont_anagrammes('tatayoyo', 'chaudchocolat'))

```

Exercice 7. On dispose dans le répertoire Ressources d'une liste de tous les mots français (selon GNU aspell¹), contenu dans un fichier nommé `listemotsfrançais.txt`. Ce fichier contient 629 570 mots : **ne pas tenter de lire directement ce fichier, sa taille est trop importante**. La structure du contenu de ce fichier est très simple : un mot par ligne et chaque ligne est terminée par `\n`.

Copier dans votre répertoire personnel le fichier `listemotsfrançais.txt` présent dans le répertoire Ressources. Avec python : écrire les instructions pour ouvrir ce fichier, récupérer le contenu du fichier, fermer le fichier. En vous appuyant sur le TP 04, créer une liste contenant la liste des mots en français. **Ne pas tenter d'afficher directement cette liste en entier, sa taille est trop importante**.

On pourra vérifier que le code fait ce qui est attendu en affichant les 10 premiers éléments de la liste, ce qui devrait donner :

```
['ah', 'an', 'ans', 'as', 'en', 'ha', 'AAAI', 'enharnacher', 'enhardi', 'enhardis']
```

Solution 7.

```

1 fichier=open('listemotsfrançais.txt','r')
2 contenu=fichier.read()
3 fichier.close()
4 listemotsfrançais = contenu.split('\n')
5 print(listemotsfrançais[0:10])

```

1. <http://aspell.net/>.

Exercice 8. Écrire une fonction `anagrammes` qui prend en argument un mot et renvoie un dictionnaire dont la clé est le mot et l'élément est la liste de toutes les anagrammes françaises de ce mot.

Attention, la liste de tous les mots français contient 629 570 mots : sur les ordinateurs de la salle de travaux pratiques, une recherche prend plus de trente secondes. Pour tester le code, il est donc recommandé d'utiliser une beaucoup plus petite liste, par exemple une liste des 50 000 premiers mots français selon GNU aspell. Avec une telle liste, l'instruction `print(anagrammes(chien))` devrait renvoyer :

```
{'chien': ['chiné']}
```

Une fois que le code est correct, afficher toutes les anagrammes des mots de votre choix en utilisant la liste complète des mots français selon GNU aspell.

Solution 8.

```

1 def anagrammes(mot):
2     '''renvoie un dictionnaire contenant un unique couple
3     dont la clé est mot et l'élément est la listes des anagrammes
4     de mot'''
5     anagrammes = {}
6     liste_anagrammes = []
7     motbis=enleve_diacritiques_majuscules_tirets(mot)
8     for motcandidat in listemotsfrançais:
9         motcandidatbis = enleve_diacritiques_majuscules_tirets(motcandidat)
10        if mot != motcandidat:
11            if sont_anagrammes(motbis,motcandidatbis):
12                liste_anagrammes.append(motcandidat)
13        anagrammes[mot] = liste_anagrammes
14    return anagrammes
15 print(anagrammes('chien'))

```

Annexe

Fonctions et méthodes sur les dictionnaires

| Opération | Exemple |
|--|--------------------------------|
| Création d'un dictionnaire vide | <code>d={}</code> |
| Création d'un dictionnaire avec un couple <code>cle/element</code> | <code>d={cle : element}</code> |
| Test d'appartenance d'une clé | <code>cle in d</code> |
| Ajout d'un couple <code>cle/element</code> | <code>d[cle]=element</code> |
| Élément correspondant à une clé | <code>d[cle]</code> |
| Test d'égalité de deux dictionnaires | <code>d1 == d2</code> |
| Ensemble des clés | <code>d.keys()</code> |
| Ensemble des éléments | <code>d.items()</code> |