

Architecture matérielle et initiation à l'algorithmique Informatique

Chapitre 1 – 06 Chaines de caractères

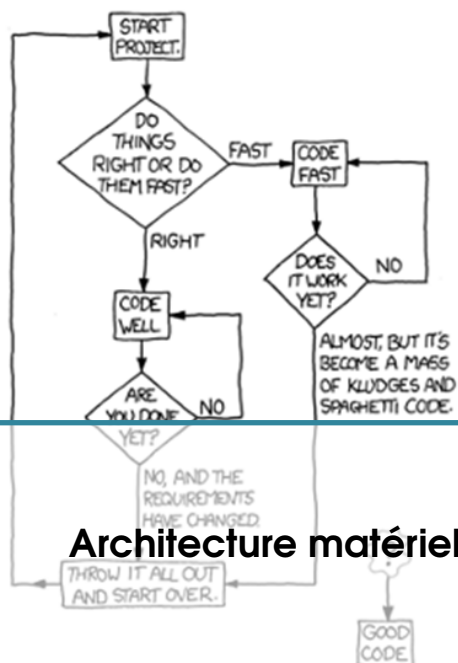
13 Novembre 2019

Savoirs et compétences :

- AA.C9 : Choisir un type de données en fonction d'un problème à résoudre
- AA.S11 : Manipulation de quelques structures de données.

1	Chaînes et type str en Python	2
1.1	Définition	2
1.2	Création et lecture	2
1.3	Opérations sur les chaînes	3
1.4	Chaînes et méthodes	3
2	Algorithme : recherche d'un mot dans un texte	4
3	Applications	5
3.1	Recherche de mots	5
3.2	Nombre d'occurrences	5

HOW TO WRITE GOOD CODE:



1 Chaînes et type `str` en Python

1.1 Définition

Pour représenter des textes, on utilise la structure de données de « chaîne de caractères » (string en anglais). En Python, cela correspond aux objets de type `str`.

Ces objets sont non-mutables, c'est-à-dire qu'une fois créés, on ne peut pas les modifier. Comme il ne peut y avoir des problèmes d'alias (par exemple, comme pour les listes Python), nous ne détaillerons pas ici leur représentation en mémoire.


1.2 Création et lecture

On définit une chaîne de caractère en entourant ces caractères par des guillemets simples, doubles, ou trois guillemets simples ou doubles. L'utilisation de guillemets simples permet d'utiliser des guillemets doubles dans la chaîne de caractère, et vice-versa.

```
>>> 'Hallo Welt'
'Hallo Welt'
>>> 'It is "Hello World" !'
'It is "Hello World" !'
>>> "Ja, aber 'not in germany' ..."
'Ja, aber "not in germany" ...'
```

On peut signaler une tabulation par `\t` et créer une nouvelle ligne par `\n`.

```
>>> s = u"\t Il faut une science politique nouvelle à un monde tout nouveau.\n[...]"
>>> print(s)
    Il faut une science politique nouvelle à un monde tout nouveau.
[...]
```

 La présence de `u` devant les guillemets permet de s'assurer de prendre les comptes les caractères spéciaux liés à l'encodage utf-8.

On accède alors aux caractères d'une chaîne en donnant son indice, comme pour les listes.

```
>>> s[0]
u'\t'
>>> s[-1]
u']'
```

Remarquons que le tranchage fonctionne similairement aux listes.

```
>>> s[14:21]
u'science'
>>> s[:11]
u'\t Il faut u'
>>> s[37:]
u'lle \xe0 un monde tout nouveau.\n[...]'
```

On peut aussi créer une chaîne de caractères à partir de nombres ou de booléens avec la fonction `str`.

```
>>> str(42)
'42'
>>> str(-3.5)
'-3.5'
>>> str(4==4.)
'True'
```

Enfin, rappelons que ce type est non mutable.

```
>>> s[0] = "a"
Traceback (most recent call last):
  File "<pyshell#1>", line 1, in <module>
    s[0] = "a"
TypeError: 'str' object does not support item assignment
```

1. Alexis de Tocqueville, De la démocratie en Amérique (1835-1840)

1.3 Opérations sur les chaînes

Comme pour les listes, on peut concaténer deux chaînes de caractères avec l'opérateur +.

```
>>> 'Hello'+'World'
'HelloWorld'
```

On peut aussi dupliquer une chaîne de caractères plusieurs fois avec l'opérateur *.

```
>>> print('GA-BU-ZO-MEU\n'*5)
GA-BU-ZO-MEU
GA-BU-ZO-MEU
GA-BU-ZO-MEU
GA-BU-ZO-MEU
GA-BU-ZO-MEU
```

La fonction `len` permet là encore de calculer la longueur d'une liste.

```
>>> len(s)
71
```

On peut aussi chercher la présence d'un motif dans une chaîne avec `in`, son absence avec `not in`

```
>>> 'Science' in s
False
>>> 'science' not in s
False
```

1.4 Chaînes et méthodes

De nombreuses méthodes existent sur les chaînes, voici les plus utiles.

Méthode	Description
<code>count(m)</code>	Compte le nombre d'occurrences du motif <code>m</code> , sans chevauchement.
<code>find(m)</code>	Renvoie la première occurrence du motif <code>m</code> .
<code>islower()</code> (et autres)	Renvoie un booléen indiquant si la chaîne est en minuscules.
<code>join(bout)</code>	Concatène les éléments de la liste <code>bout</code> , séparés par la chaîne.
<code>lower()</code>	Met en minuscule la chaîne de caractères.
<code>split(sep)</code>	Sépare la chaîne selon le séparateur <code>sep</code>
<code>strip(char)</code>	Enlève les lettres en début/fin si elles sont dans la chaîne <code>char</code> .
<code>upper()</code>	Mets en majuscules la chaîne.

FIGURE 1 – Quelques méthodes sur les chaînes de caractères.

```
>>> s = "    que ; d'espaces ; mes ; amis ! ;    "
>>> s.count('es')
3
>>> s.find('es')
12
>>> s.islower()
True
>>> '+'.join([str(i) for i in range(14)])
'0+1+2+3+4+5+6+7+8+9+10+11+12+13'
>>> 'HAHAHAHAHA !'.lower()
'hahahahaha !'
>>> s.split(';')
['    que ', ' " d'espaces ', ' mes ', ' amis ! ', '    ']
>>> s.split()
['que', ' ', ' ', ' "d'espaces", ', ' ', 'mes', ' ', ' ', 'amis', '!', ' ', ' ;']
>>> s.strip()
"que ; d'espaces ; mes ; amis ! ;"
>>> s.strip('q; ')
"ue ; d'espaces ; mes ; amis !"
>>> s.upper()
"    QUE ; D'ESPACES ; MES ; AMIS ! ;    "
```

2 Algorithme : recherche d'un mot dans un texte

Le seul algorithme sur les chaînes que vous devez pouvoir reproduire est celui de la recherche d'un mot m dans une chaîne de caractères s .

Nous allons écrire un algorithme naïf. On commence par calculer la longueur de m . Ensuite, on parcourt séquentiellement la chaîne s en regardant à chaque fois si le mot se trouve au début de l'emplacement considéré.

Cela donne l'écriture suivante.

```
def recherche(m,s):
    """Recherche le mot m dans la chaîne s
    Préconditions : m et s sont des chaînes de caractères"""
    long_s = len(s) # Longueur de s
    long_m = len(m) # Longueur de m
    for i in range(long_s-long_m+1):
        # Invariant : m n'a pas été trouvé dans s[0:i+long_m-1]
        if s[i:i+long_m] == m: # On a trouvé m
            return True
    return False

print(recherche('42','Les 7 samourais'))
print(recherche('7','Les 7 samourais'))
```

Cela renvoie alors :

```
False
True
```



- Remarquons que l'on aurait pu exécuter cet algorithme sur une liste Python...
- Nous pouvons aussi remarquer que l'algorithme précédent utilise une tranche, qui est particulière à Python, et que nous ne savons pas vraiment comment se comporte la comparaison de tableaux. On peut ainsi choisir d'explicitement complètement cette partie là comme suit.

```
def recherche(m,s):
    """Recherche le mot m dans la chaîne s
    Préconditions : m et s sont des chaînes de caractères"""
    long_s = len(s) # Longueur de s
    long_m = len(m) # Longueur de m
    for i in range(long_s-long_m+1):
        # Invariant : m n'a pas été trouvé dans s[0:i+long_m-1]
        j = 0
        while j < long_m and m[j] == s[i+j]:
            # Invariant : m[:j] == s[i:i+j]
            j = j+1
            # Invariant : m[:j] == s[i:i+j]
        if j == long_m:
            # Invariant précédent : m == s[i:i+long_m]
            return True
    return False

print(recherche('42','Les 7 samourais'))
print(recherche('sam','Les 7 samourais'))
```

On vérifie, cela renvoie bien :

```
False
True
```

TD

Savoirs et compétences :

- ☐ AA.C9 : Choisir un type de données en fonction d'un problème à résoudre
- ☐ AA.S11 : Manipulation de quelques structures de données.

3 Applications

3.1 Recherche de mots

Q 1 : Écrire une fonction `c_recherche(m, s)` de recherche indiquant si le mot `m` est dans la chaîne `s` insensible à la casse (majuscules ou minuscules).

Q 2 : Écrire une fonction `i_recherche(m, s)` de recherche du mot `m` dans une chaîne `s`, renvoyant la position de la première occurrence du mot dans la chaîne et une erreur si le mot n'est pas dans le texte.

3.2 Nombre d'occurrences

Q 1 : Écrire une fonction `compte(m, s)` comptant le nombre d'occurrences du mot `m` dans une chaîne `s`, chevauchements compris.

Q 2 : Écrire une fonction `dist_compte(m, s)` comptant le nombre d'occurrences du mot `m` dans une chaîne `s`, sans chevauchements.