

Chapitre 11 : Les chaînes de caractères

I. Introduction :

Une chaîne de caractères est une suite ordonnée de caractères.

En algorithmique, une chaîne de caractères est considérée comme un tableau de caractères.

En Python, les chaînes de caractères sont classées dans les séquences tout comme les listes et les tuples (**les tuples et les chaînes sont des séquences immuables, c.à.d. qu'elles ne peuvent pas être modifiées**).

Et on a vu que pour déclarer une chaîne on doit utiliser les guillemets ou les simples quotes. On a vu aussi comment on peut les manipuler en les concaténant ensemble ou les parcourir avec la boucle for, mais en fait avec les chaînes (string) on n'arrête pas ici mais il y a plein de méthodes avec lesquelles on peut manipuler les chaînes

II. Code ASCII :

Dans un système informatique, à chaque caractère est associé une valeur numérique : son code ASCII (American Standard Code for Information Interchange). L'ensemble des codes est recensé dans une table nommée "table des codes ASCII". Quand on stocke un caractère en mémoire (dans une variable), on mémorise en réalité son code ASCII. Un code ASCII est codé sur un octet (huit bits).

- **Comparaison des caractères :**

Voici un exemple de comparaison de caractères :

"B" > "A" car le code ASCII du caractère "B" (66 en base 10) est supérieur au code ASCII du caractère "A" (65 en base 10).

- Pour comparer deux chaînes de caractères, on compare les caractères de même rang dans les deux chaînes en commençant par le premier caractère de chaque chaîne (le premier caractère de la première chaîne

est comparé au premier caractère de la seconde chaîne, le deuxième caractère de la première chaîne est comparé au deuxième caractère de la seconde chaîne, et ainsi de suite...).

Exemples : Comparaison de deux chaînes

- **"baobab" < "sapin"** car le code ASCII de "b" est inférieur au code ASCII de "s".
- **"baobab" > "banania"** car le code ASCII de "o" est supérieur au code ASCII de "n" (la comparaison ne peut pas se faire sur les deux premiers caractères car ils sont identiques).
- **"1999" > "1998"** car le code ASCII de "9" est supérieur au code ASCII de "8" (la comparaison ne peut pas se faire sur les trois premiers caractères car ils sont identiques). Attention, ici ce ne sont pas des valeurs numériques qui sont comparées, mais bien des caractères.
- **"333" > "1230"** car le code ASCII de "3" est supérieur au code ASCII de "1".
- **"333" < "3330"** car la seconde chaîne a une longueur supérieure à celle de la première (la comparaison ne peut pas se faire sur les trois premiers caractères car ils sont identiques).
- **"Baobab" < "baobab"** car le code ASCII de "b" est supérieur au code ASCII de "B".

III. La manipulation des chaînes de caractère en algorithmes :

- La fonction **SSCHAINE (chaîne , position , nombre) : chaîne**

Le rôle de la fonction SSCHAINE() est de retourner une "sous-chaîne" (copie d'un extrait de la chaîne passée en premier paramètre) de nombre caractères de la chaîne chaîne à partir du caractère qui se trouve en position **position**

- La longueur de la chaîne : la fonction **LONGUEUR(chaîne) :entier**

La fonction `LONGUEUR()` a pour rôle de calculer et de retourner le nombre de caractères présents dans la chaîne de caractères passée en paramètre.

- La fonction **`RANG(chaine , sous-chaine , position)` : entier**

La fonction `RANG()` recherche et retourne la première occurrence de la sous-chaîne dans la chaîne. La recherche commence à partir de position. Cette fonction retourne la position du premier caractère de la sous-chaîne, ou la valeur -1 si la "sous-chaîne" n'existe pas dans la chaîne (à partir de position).

- La fonction **`CODE(caractere)` : entier**

La fonction `CODE()` retourne le code ASCII du caractère passé en paramètre.

- La fonction **`CAR(code_ascii)` : caractere**

La fonction `CAR()` retourne le caractère dont le code ASCII est passé en paramètre.

- La fonction **`MINUSCULE(chaine)` : chaine**

Retourne la même chaine passée en paramètre en minuscule

- La fonction **`MAJUSCULE(chaine)` : chaine**

Retourne la même chaine passée en paramètre en majuscule

IV. La manipulation des chaines de caractère en python :

Les chaînes de caractères peuvent être considérées comme des listes (de caractères) un peu particulières :

```
1 | >>> animaux = "girafe tigre"
2 | >>> animaux
3 | 'girafe tigre'
4 | >>> len(animaux)
5 | 12
6 | >>> animaux[3]
7 | 'a'
```

Nous pouvons donc utiliser certaines propriétés des listes comme les tranches :

```

1  >>> animaux = "girafe tigre"
2  >>> animaux[0:4]
3  'gira'
4  >>> animaux[9:]
5  'gre'
6  >>> animaux[:-2]
7  'girafe tig'
8  >>> animaux[1:-2:2]
9  'iaetg'

```

Mais contrairement aux listes, les chaînes de caractères présentent toutefois une différence notable, ce sont **des listes non modifiables**. Une fois une chaîne de caractères définie, vous ne pouvez plus modifier un de ses éléments. Le cas échéant, Python renvoie un message d'erreur :

```

1  >>> animaux = "girafe tigre"
2  >>> animaux[4]
3  'f'
4  >>> animaux[4] = "F"
5  Traceback (most recent call last):
6    File "<stdin>", line 1, in <module>
7  TypeError: 'str' object does not support item assignment

```

Par conséquent, si vous voulez modifier une chaîne de caractères, vous devez en construire une nouvelle. Pour cela, n'oubliez pas que les opérateurs de concaténation (+) et de duplication (*) (introduits dans le chapitre 2 *Variables*) peuvent vous aider. Vous pouvez également générer une liste, qui elle est modifiable, puis revenir à une chaîne de caractères (voir plus bas).

V. Caractères spéciaux :

Il existe certains caractères spéciaux comme `\n` que nous avons déjà vu (pour le retour à la ligne). Le caractère `\t` produit une tabulation. Si vous voulez écrire des guillemets simples ou doubles et que ceux-ci ne soient pas confondus avec les guillemets de déclaration de la chaîne de caractères, vous pouvez utiliser `\'` ou `\"`.

```

1  >>> print("Un retour à la ligne\npuis une tabulation\t puis un guillemet'")
2  Un retour à la ligne
3  puis une tabulation      puis un guillemet"
4  >>> print('J\'affiche un guillemet simple')
5  J'affiche un guillemet simple

```

Vous pouvez aussi utiliser astucieusement des guillemets doubles ou simples pour déclarer votre chaîne de caractères :

```
1 >>> print("Un brin d'ADN")
2 Un brin d'ADN
3 >>> print('Python est un "super" langage de programmation')
4 Python est un "super" langage de programmation
```

Quand on souhaite écrire un texte sur plusieurs lignes, il est très commode d'utiliser les guillemets triples qui conservent le formatage (notamment les retours à la ligne) :

```
1 >>> x = """souris
2 ... chat
3 ... abeille"""
4 >>> x
5 'souris\nchat\nabeille'
6 >>> print(x)
7 souris
8 chat
9 abeille
```

VI. Méthodes associées aux chaînes de caractères

Voici quelques [méthodes](#) spécifiques aux objets de type str :

```
1 >>> x = "girafe"
2 >>> x.upper()
3 'GIRAFE'
4 >>> x
5 'girafe'
6 >>> 'TIGRE'.lower()
7 'tigre'
```

- Les méthodes **.lower()** et **.upper()** renvoient un texte en minuscule et en majuscule respectivement. On remarque que l'utilisation de ces méthodes n'altère pas la chaîne de caractères de départ mais renvoie une chaîne de caractères transformée.

Pour mettre en majuscule la première lettre seulement, vous pouvez faire :

```
1 | >>> x[0].upper() + x[1:]
2 | 'Girafe'
```

- Les méthodes **.isupper()** et **.islower()** retournent True si la chaîne est toute en majuscule ou minuscule
- La méthode **.capitalize()** plus simplement utiliser la méthode adéquate :

```
1 | >>> x.capitalize()
2 | 'Girafe'
```

- la méthode **.split()** Il existe une méthode associée aux chaînes de caractères qui est particulièrement pratique :

```
1 | >>> animaux = "girafe tigre singe souris"
2 | >>> animaux.split()
3 | ['girafe', 'tigre', 'singe', 'souris']
4 | >>> for animal in animaux.split():
5 | ...     print(animal)
6 | ...
7 | girafe
8 | tigre
9 | singe
10 | souris
```

La méthode **.split()** découpe une chaîne de caractères en plusieurs éléments appelés *champs*, en utilisant comme séparateur n'importe quelle combinaison « d'espace(s) blanc(s) ».

Un espace blanc (*whitespace* en anglais) correspond aux caractères qui sont invisibles à l'œil, mais qui occupent de l'espace dans un texte. Les espaces blancs les plus classiques sont l'espace, la tabulation et le retour à la ligne.

Il est possible de modifier le séparateur de champs, par exemple :

```
1 | >>> animaux = "girafe:tigre:singe::souris"
2 | >>> animaux.split(":")
3 | ['girafe', 'tigre', 'singe', '', 'souris']
```

Attention, dans cet exemple, le séparateur est un seul caractère « : » (et non pas une combinaison de un ou plusieurs :) conduisant ainsi à une chaîne vide entre singe et souris.

Il est également intéressant d'indiquer à `.split()` le nombre de fois qu'on souhaite découper la chaîne de caractères avec l'argument `maxsplit` :

```
1 | >>> animaux = "girafe tigre singe souris"
2 | >>> animaux.split(maxsplit=1)
3 | ['girafe', 'tigre singe souris']
4 | >>> animaux.split(maxsplit=2)
5 | ['girafe', 'tigre', 'singe souris']
```

- La méthode **`.find(sous-chaine , position)`**, quant à elle, recherche une chaîne de caractères passée en argument :

```
1 | >>> animal = "girafe"
2 | >>> animal.find("i")
3 | 1
4 | >>> animal.find("afe")
5 | 3
6 | >>> animal.find("z")
7 | -1
8 | >>> animal.find("tig")
9 | -1
```

Si l'élément recherché est trouvé, alors l'indice du début de l'élément dans la chaîne de caractères est renvoyé. Si l'élément n'est pas trouvé, alors la valeur -1 est renvoyée.

Si l'élément recherché est trouvé plusieurs fois, seul l'indice de la première occurrence est renvoyé :

```
1 | >>> animaux = "girafe tigre"
2 | >>> animaux.find("i")
3 | 1
```

- La méthode **`.index(sous-chaine , position)`** : recherche une sous chaîne et renvoie sa position si la sous chaîne est trouvée sinon génère une exception (`ValueError`)

- La fonction **ord(caractere)** : renvoie le code ascii d'un caractère passé en paramètre .
- La fonction **chr(code_ascii)** : renvoie le caractère correspondant au code ascii passé en paramètre
- la méthode **.replace(ancienne_chaine , nouvelle_chaine)** qui substitue une chaîne de caractères par une autre :

```

1 | >>> animaux = "girafe tigre"
2 | >>> animaux.replace("tigre", "singe")
3 | 'girafe singe'
4 | >>> animaux.replace("i", "o")
5 | 'gorafe togre'

```

- La méthode **.count(chaine)** compte le nombre d'occurrences d'une chaîne de caractères passée en argument :

```

1 | >>> animaux = "girafe tigre"
2 | >>> animaux.count("i")
3 | 2
4 | >>> animaux.count("z")
5 | 0
6 | >>> animaux.count("tigre")
7 | 1

```

- La méthode **.startswith(), .endswith()** vérifie si une chaîne de caractères commence par une autre chaîne de caractères :

```

1 | >>> chaine = "Bonjour monsieur le capitaine !"
2 | >>> chaine.startswith("Bonjour")
3 | True
4 | >>> chaine.startswith("Au revoir")
5 | False

```

Cette méthode est particulièrement utile lorsqu'on lit un fichier et que l'on veut récupérer certaines lignes commençant par un mot-clé. Par exemple dans un fichier PDB, les lignes contenant les coordonnées des atomes commencent par le mot-clé ATOM.

- la méthode **.strip()** permet de « nettoyer les bords » d'une chaîne de caractères :


```
1 | >>> chaine = " Comment enlever les espaces au début et à la fin ?  
2 | >>> chaine.strip()  
3 | 'Comment enlever les espaces au début et à la fin ?'
```

La méthode **.strip()** enlève les espaces situés sur les bords de la chaîne de caractère mais pas ceux situés entre des caractères visibles. En réalité, cette méthode enlève n'importe quel combinaison « d'espace(s) blanc(s) » sur les bords, par exemple :

```
1 | >>> chaine = " \tfonctionne avec les tabulations et les retours à la ligne\n"  
2 | >>> chaine.strip()  
3 | 'fonctionne avec les tabulations et les retours à la ligne'
```

La méthode **.strip()** est très pratique quand on lit un fichier et qu'on veut se débarrasser des retours à la ligne.

VII. Extraction de valeurs numériques d'une chaîne de caractères :

Une tâche courante en Python est de lire une chaîne de caractères (provenant par exemple d'un fichier), d'extraire des valeurs de cette chaîne de caractères puis ensuite de les manipuler.

On considère par exemple la chaîne de caractères val :

```
1 | >>> val = "3.4 17.2 atom"
```

On souhaite extraire les valeurs 3.4 et 17.2 pour ensuite les additionner.

Dans un premier temps, on découpe la chaîne de caractères avec la méthode **.split()** :

```
1 | >>> val2 = val.split()  
2 | >>> val2  
3 | ['3.4', '17.2', 'atom']
```

On obtient alors une liste de chaînes de caractères. On transforme ensuite les deux premiers éléments de cette liste en *floats* (avec la fonction **float()**) pour pouvoir les additionner :

```
1 | >>> float(val2[0]) + float(val2[1])
2 | 20.599999999999998
```

Remarque : Retenez bien l'utilisation des instructions précédentes pour extraire des valeurs numériques d'une chaîne de caractères. Elles sont régulièrement employées pour analyser des données depuis un fichier.

VIII. Conversion d'une liste de chaînes de caractères en une chaîne de caractères :

On a vu dans le chapitre 2 *Variables* la conversion d'un type simple (entier, *float* et chaîne de caractères) en un autre avec les fonctions `int()`, `float()` et `str()`. La conversion d'une liste de chaînes de caractères en une chaîne de caractères est un peu particulière puisqu'elle fait appel à la méthode `.join()`.

```
1 | >>> seq = ["A", "T", "G", "A", "T"]
2 | >>> seq
3 | ['A', 'T', 'G', 'A', 'T']
4 | >>> "-".join(seq)
5 | 'A-T-G-A-T'
6 | >>> " ".join(seq)
7 | 'A T G A T'
8 | >>> "".join(seq)
9 | 'ATGAT'
```

Les éléments de la liste initiale sont concaténés les uns à la suite des autres et intercalés par un séparateur qui peut être n'importe quelle chaîne de caractères. Ici, on a utilisé un tiret, un espace et rien (une chaîne de caractères vide).

Attention, la méthode `.join()` ne s'applique qu'à une liste de chaînes de caractères.

```
1 | >>> maliste = ["A", 5, "G"]
2 | >>> " ".join(maliste)
3 | Traceback (most recent call last):
4 |   File "<stdin>", line 1, in <module>
5 | TypeError: sequence item 1: expected string, int found
```

On espère qu'après ce petit tour d'horizon vous serez convaincu de la richesse des méthodes associées aux chaînes de caractères. Pour avoir une liste exhaustive de l'ensemble des méthodes associées à une variable particulière, vous pouvez utiliser la fonction `dir()`.

```
1 >>> animaux = "girafe tigre"
2 >>> dir(animaux)
3 ['__add__', '__class__', '__contains__', '__delattr__', '__dir__',
4  '__doc__', '__eq__', '__format__', '__ge__', '__getattribute__', '_
5  _getitem__', '__getnewargs__', '__gt__', '__hash__', '__init__', '_
6  _init_subclass__', '__iter__', '__le__', '__len__', '__lt__', '__mo
7  d__', '__mul__', '__ne__', '__new__', '__reduce__', '__reduce_ex__'
8  , '__repr__', '__rmod__', '__rmul__', '__setattr__', '__sizeof__',
9  '__str__', '__subclasshook__', 'capitalize', 'casefold', 'center',
10 'count', 'encode', 'endswith', 'expandtabs', 'find', 'format', 'for
11 mat_map', 'index', 'isalnum', 'isalpha', 'isdecimal', 'isdigit', 'i
12 sidentifier', 'islower', 'isnumeric', 'isprintable', 'isspace', 'is
13 title', 'isupper', 'join', 'ljust', 'lower', 'lstrip', 'maketrans',
14 'partition', 'replace', 'rfind', 'rindex', 'rjust', 'rpartition',
15 'rsplit', 'rstrip', 'split', 'splitlines', 'startswith', 'strip',
16 'swapcase', 'title', 'translate', 'upper', 'zfill']
```

Pour l'instant, vous pouvez ignorer les méthodes qui commencent et qui se terminent par deux tirets bas (*underscores*) `__`.

Vous pouvez également accéder à l'aide et à la documentation d'une méthode particulière avec `help()`, par exemple pour la méthode `.split()` :

```
1 >>> help(animaux.split)
2 Help on built-in function split:
3
4 split(...)
5     S.split([sep [,maxsplit]]) -> list of strings
6
7     Return a list of the words in the string S, using sep as the
8     delimiter string.  If maxsplit is given, at most maxsplit
9     splits are done.  If sep is not specified or is None, any
10    whitespace string is a separator.
11 (END)
```

Attention à ne pas mettre les parenthèses à la suite du nom de la méthode. L'instruction correcte est `help(animaux.split)` et non pas `help(animaux.split())`

IX. Exercices :

Exercice 1 : Ecrivez un algorithme puis un programme en python qui permet de convertir une chaîne saisie par l'utilisateur en majuscule

Exercice 2 : Ecrivez un algorithme puis un programme en python qui compte le nombre d'occurrences d'une lettre donnée dans cette chaîne.

Exercice 3 : Ecrivez un algorithme puis un programme en python qui compte le nombre d'occurrences des lettres et nombres.

Exercice 4 : Ecrivez un algorithme puis un programme en python qui permet de calculer le nombre de répétition de chaque caractère de la chaîne.

Exercice 5 : Ecrivez un algorithme puis un programme en python qui permet de calculer le nombre de voyelles dans une chaîne de caractères.

Exercice 6 : Ecrivez un algorithme puis un programme en python qui supprime toutes les lettres « e » (minuscules) d'un texte de moins d'une ligne (supposée ne pas dépasser 132 caractères) fourni au clavier.

Exercice 7 : Ecrivez un algorithme puis un programme en python qui supprime toutes les voyelles d'un texte de moins d'une ligne fourni au clavier.

Exercice 8 : Ecrivez un algorithme puis un programme en python qui permet de supprimer les doublons dans une chaîne de caractère saisie par l'utilisateur.

Exercice 9 : Soit la liste ['girafe', 'tigre', 'singe', 'souris'].

1. Affichez chaque élément ainsi que sa taille (nombre de caractères).
2. Affichez pour chaque élément le nombre des voyelles et le nombre des consonnes.