

Cours de Python



<https://python.sdv.univ-paris-diderot.fr/>

Patrick Fuchs et Pierre Poulain

prénom [point] nom [arobase] univ-paris-diderot [point] fr

version du 7 septembre 2018

Université Paris Diderot-Paris 7, Paris, France

Table des matières

1	Introduction	11
1.1	Quelques mots sur l'origine de ce cours	11
1.2	Remerciements	11
1.3	Avant de commencer	12
1.4	C'est quoi Python ?	12
1.5	Premier contact avec Python sous Linux	13
1.6	Premier programme Python	14
1.6.1	Appel de l'interpréteur	14
1.6.2	Appel direct du script	14
1.7	Commentaires	15
1.8	Notion de bloc d'instructions et d'indentation	15
1.9	Python 2 ou Python 3 ?	16
2	Variables	17
2.1	Définition d'une variable	17
2.2	Les types de variables	18
2.3	Nommage des variables	18
2.4	Opérations	18
2.4.1	Opérations sur les types numériques	18
2.4.2	Opérations sur les chaînes de caractères	19
2.4.3	Opérations illicites	19
2.5	La fonction type()	20
2.6	Conversion de types	20
2.7	Note sur la division	20
2.8	Note sur le vocabulaire et la syntaxe	21
3	Affichage	23
3.1	Écriture formatée	23
3.2	Ancienne méthode de formatage des chaînes de caractères	26
3.3	Note sur le vocabulaire et la syntaxe	27
3.4	Exercices	27
3.4.1	Affichage dans l'interpréteur et dans un programme	27
3.4.2	Poly-A	27
3.4.3	Poly-A et poly-GC	27
3.4.4	Écriture formatée	27
3.4.5	Écriture formatée 2	28
4	Listes	29
4.1	Définition	29
4.2	Utilisation	29
4.3	Opération sur les listes	30
4.4	Indiçage négatif	31

4.5	Tranches	31
4.6	Fonction len()	32
4.7	Les fonctions range() et list()	32
4.8	Listes de listes	33
4.9	Exercices	33
4.9.1	Jours de la semaine	34
4.9.2	Saisons	34
4.9.3	Table des 9	34
4.9.4	Nombres pairs	34
5	Boucles et comparaisons	35
5.1	Boucles for	35
5.1.1	Principe	35
5.1.2	Fonction range()	37
5.1.3	Nommage de la variable d'itération	37
5.1.4	Itération sur les indices	37
5.2	Comparaisons	39
5.3	Boucles while	40
5.4	Exercices	40
5.4.1	Boucles de base	41
5.4.2	Boucle et jours de la semaine	41
5.4.3	Table des 1	41
5.4.4	Nombres pairs et impairs	41
5.4.5	Calcul de moyenne	41
5.4.6	Produit de nombres consécutifs	41
5.4.7	Triangle	41
5.4.8	Triangle inversé	42
5.4.9	Triangle gauche	42
5.4.10	Triangle isocèle	42
5.4.11	Parcours de matrice	43
5.4.12	Parcours de demi-matrice sans la diagonale (exercice ++)	43
5.4.13	Sauts de puce	44
5.4.14	Suite de Fibonacci (exercice +++)	44
6	Tests	47
6.1	Définition	47
6.2	Tests à plusieurs cas	47
6.3	Tests multiples	49
6.4	Instructions break et continue	50
6.5	Tests de valeur sur des réels	50
6.6	Exercices	51
6.6.1	Jours de la semaine	51
6.6.2	Séquence complémentaire d'un brin d'ADN	51
6.6.3	Minimum d'une liste	51
6.6.4	Fréquence des acides aminés	51
6.6.5	Notes et mention d'un étudiant	52
6.6.6	Nombres pairs	52
6.6.7	L'énigme du père Fouras	52
6.6.8	Conjecture de Syracuse (exercice +++)	52
6.6.9	Attribution simple de la structure secondaire des résidus d'une protéine (exercice +++)	53
6.6.10	Détermination des nombres premiers inférieurs à 100 (exercice +++)	53
6.6.11	Recherche d'un nombre par dichotomie (exercice +++)	54

7 Fichiers	57
7.1 Lecture dans un fichier	57
7.1.1 Méthode read()	58
7.1.2 Méthode readline()	58
7.1.3 Méthodes seek() et tell()	59
7.1.4 Itérations directement sur le fichier	59
7.2 Écriture dans un fichier	60
7.3 Méthode optimisée d'ouverture et de fermeture de fichier	60
7.4 Note sur les retours chariots sous Unix et sous Windows	61
7.5 Importance des conversions de types avec les fichiers	62
7.6 Exercices	62
7.6.1 Lecture et saut de ligne	62
7.6.2 Écriture et saut de ligne	62
7.6.3 Structures secondaires	62
7.6.4 Spirale (exercice +++)	62
8 Modules	65
8.1 Définition	65
8.2 Importation de modules	65
8.3 Obtenir de l'aide sur les modules importés	66
8.4 Modules courants	67
8.5 Module sys : passage d'arguments	68
8.6 Module os	69
8.7 Exercices	70
8.7.1 Racine carrée	70
8.7.2 Cosinus	70
8.7.3 Liste de fichiers dans un répertoire	70
8.7.4 Affichage temporisé	70
8.7.5 Séquences aléatoires de chiffres	70
8.7.6 Séquences aléatoires de bases	71
8.7.7 Jour de naissance	71
8.7.8 Détermination du nombre pi par la méthode Monte Carlo (exercice +++)	71
9 Fonctions	73
9.1 Principe et généralités	73
9.2 Définition	74
9.3 Passage d'arguments	75
9.4 Variables locales et variables globales	76
9.5 Exercices	80
9.5.1 Fonctions et pythontutor	80
9.5.2 Fonction puissance	81
9.5.3 Fonction pyramide	81
9.5.4 Fonction nombre premier	81
9.5.5 Fonction complement	81
9.5.6 Fonction distance	82
9.5.7 Fonctions distribution et stat	82
9.5.8 Fonction distance à l'origine	82
9.5.9 Fonction aire sous la courbe (exercice +++)	82
10 Plus sur les chaînes de caractères	85
10.1 Préambule	85
10.2 Chaînes de caractères et listes	85
10.3 Caractères spéciaux	86

10.4	Méthodes associées aux chaînes de caractères	86
10.5	Extraction de valeurs numériques d'une chaîne de caractères	88
10.6	Conversion d'une liste de chaînes de caractères en une chaîne de caractères	88
10.7	Exercices	90
10.7.1	Parcours d'une liste de chaînes de caractères	90
10.7.2	Fréquence des bases dans une séquence nucléique	90
10.7.3	Conversion des acides aminés du code à trois lettres au code à une lettre	90
10.7.4	Distance de Hamming	91
10.7.5	Palindrome	91
10.7.6	Mot composable	91
10.7.7	Alphabet et pangramme	91
10.7.8	Affichage des carbones alpha d'une structure de protéine	92
10.7.9	Calcul des distances entre les carbones alpha consécutifs d'une structure de protéine	92
11	Plus sur les listes	95
11.1	Propriétés des listes	95
11.2	Test d'appartenance	97
11.3	Copie de listes	97
11.4	Exercices	99
11.4.1	Tri de liste	99
11.4.2	Séquence nucléique aléatoire	99
11.4.3	Séquence nucléique complémentaire	99
11.4.4	Doublons	100
11.4.5	Séquence nucléique aléatoire 2	100
11.4.6	Triangle de Pascal (Exercice +++)	100
12	Plus sur les fonctions	101
12.1	Appel d'une fonction dans une fonction	101
12.2	Fonctions récursives	102
12.3	Portée des variables	103
12.4	Portée des listes	105
12.5	Règle LGI	106
12.6	Recommandations	106
12.7	Exercices	107
12.7.1	Prédire la sortie	107
13	Dictionnaires et tuples	109
13.1	Dictionnaires	109
13.1.1	Itération sur les clefs pour obtenir les valeurs	109
13.1.2	Méthodes keys() et values()	110
13.1.3	Existence d'une clef	110
13.1.4	Liste de dictionnaires	110
13.2	Tuples	111
13.3	Exercices	112
13.3.1	Composition en acides aminés	112
13.3.2	Mots de 2 lettres	112
13.3.3	Mots de 3 et 4 lettres	112
13.3.4	Mots de 2 lettres de <i>Saccharomyces cerevisiae</i>	112
13.3.5	Mots de n lettres et fichiers genbank	112
13.3.6	Mots de n lettres du génome d' <i>E. Coli</i>	112
13.3.7	Dictionnaire et carbone alpha	113
13.3.8	Dictionnaire et PDB	113
13.3.9	Barycentre d'une protéine	113

14 Création de modules	115
14.1 Pourquoi créer ses propres modules ?	115
14.2 Comment créer son propre module	115
14.3 Utilisation de son propre module	116
14.4 Les docstrings	116
14.5 Modules et fonctions	117
14.6 Module ou script ?	117
14.7 Exercices	119
14.7.1 Module ADN	119
15 Expressions régulières et parsing	121
15.1 Définition et syntaxe	121
15.2 Module re et fonction search	123
15.2.1 La fonction search()	123
15.2.2 Les fonctions match() et fullmatch()	124
15.2.3 Compilation d'expressions régulières	124
15.2.4 Groupes	125
15.2.5 La fonction findall()	125
15.2.6 La fonction sub()	126
15.3 Exercices	127
15.3.1 Regex de base	127
15.3.2 Nettoyeur d'espaces	127
15.3.3 Nettoyeur de doublons	127
15.3.4 Compteur de gènes	127
15.3.5 Recherche de séquence dans un fichier gbk	127
15.3.6 Le défi du dé-htmliseur (+++)	128
15.3.7 Remplacements puissants (+++)	128
16 Autres modules d'intérêt	129
17 Modules d'intérêt en bioinformatique	131
17.1 Module numpy	131
17.1.1 Objets de type array	132
17.1.2 array et dimensions	133
17.1.3 Indices	135
17.1.4 Construction automatique de matrices	135
17.1.5 Un peu d'algèbre linéaire	136
17.1.6 Un peu de transformée de Fourier	137
17.2 Module biopython	138
17.3 Module matplotlib	139
17.3.1 Représentation sous forme de points	139
17.3.2 Représentation sous forme de courbe	141
17.3.3 Représentation sous forme de barres	142
17.4 Exercices	143
17.4.1 Extraction des coordonnées atomiques	143
17.4.2 Lecture des coordonnées	144
17.4.3 Utilisation de numpy	144
17.4.4 Calcul de la distance	144
18 Jupyter et ses notebooks	145
18.1 Lancement de Jupyter et création d'un notebook	145
18.2 Le format Markdown	148
18.3 Des graphiques dans les notebooks	150

18.4	Les <i>magic commands</i>	151
19	Avoir la classe avec les objets	153
19.1	Construction d'une classe	153
19.2	Utilisation de la classe Rectangle	154
19.3	Modification du rectangle en carré	154
19.4	Création d'un objet Rectangle avec des paramètres imposés	154
19.5	Exercices	154
19.5.1	Classe Rectangle	154
19.5.2	Classe Atome	155
19.5.3	Classe Atome améliorée	155
20	Pour aller plus loin	157
20.1	Shebang et /usr/bin/env python3	157
20.2	Différences Python 2 et Python 3	157
20.2.1	Interpréteur par défaut	157
20.2.2	La fonction print()	158
20.2.3	Division d'entiers	158
20.2.4	La fonction range()	158
20.2.5	Encodage et utf-8	159
20.3	Liste de compréhension	160
20.3.1	Nombres pairs compris entre 0 et 30	160
20.3.2	Jeu sur la casse des mots d'une phrase	160
20.3.3	Formatage d'une séquence avec 60 caractères par ligne	160
20.3.4	Formatage fasta d'une séquence (avec la ligne de commentaire)	160
20.3.5	Sélection des carbones alpha dans un fichier pdb	161
20.4	Gestion des erreurs	161
20.5	Sauvegardez votre historique de commandes	163
21	Mini-projets	165
21.1	Mots anglais dans le protéome humain	165
21.1.1	Composition aminée	165
21.1.2	Des mots	165
21.1.3	Des protéines	165
21.1.4	À la pêche aux mots	166
21.1.5	Et le mot le plus fréquent est...	166
21.1.6	Pour être plus complet	166
21.2	genbank2fasta	167
21.2.1	genbank2fasta sans expression régulière	167
21.2.2	genbank2fasta avec expression régulière	170
22	Annexe : Quelques formats de données rencontrés en biologie	175
22.1	Fasta	175
22.1.1	Exemples	175
22.1.2	Manipulation avec Python	176
22.2	GenBank	177
22.2.1	L'en-tête	177
22.2.2	Les <i>features</i>	178
22.2.3	La séquence	179
22.2.4	Manipulation avec Python	179
22.3	PDB	180
22.3.1	En-tête	180
22.3.2	Coordonnées	181

22.3.3 Manipulation avec Python 184

Chapitre 1

Introduction

1.1 Quelques mots sur l'origine de ce cours

Ce cours a été conçu à l'origine pour les étudiants débutants en programmation Python des filières de biologie et de biochimie de l'Université Paris Diderot - Paris 7¹; et plus spécialement pour les étudiants du master Biologie Informatique.

Ce cours est basé sur la version 3 de Python, version recommandée par la communauté scientifique. Des références à l'ancienne version, Python 2, seront néanmoins régulièrement apportées.

Si vous relevez des erreurs à la lecture de ce document, merci de nous les signaler.

Le cours est disponible en version HTML² et PDF³.

1.2 Remerciements

Un grand merci à Sander⁴ du *Centre for Molecular and Biomolecular Informatic* de Nijmegen aux Pays-Bas pour la toute première version⁵ de ce cours qui remonte à l'année 2003.

Nous remercions le professeur Philip Guo⁶ de l'UC San Diego, pour nous avoir autorisé à utiliser des clichés de son excellent site pythontutor⁷.

Merci également à tous les contributeurs, occasionnels ou réguliers : Jennifer Becq, Virginie Martiny, Romain Laurent, Benoist Laurent, Benjamin Boyer, Hubert Santuz, Catherine Lesourd, Philippe Label, Rémi Cuchillo, Cédric Gageat, Philibert Malbranche, Mikaël Naveau, Amélie Bacle, Alexandra Moine-Franel.

Nous remercions aussi Denis Mestivier de qui nous nous sommes inspirés pour certains exercices.

Enfin, merci à vous tous, les curieux de Python, qui avez été nombreux à nous envoyer des retours sur ce cours, à nous suggérer des améliorations et à nous signaler des coquilles.

De nombreuses personnes nous ont aussi demandé les corrections des exercices. Nous ne les mettons pas sur le site afin d'éviter la tentation de les regarder trop vite, mais vous pouvez nous écrire et nous vous les enverrons.

1. <http://www.univ-paris-diderot.fr/>

2. <https://python.sdv.univ-paris-diderot.fr/index.html>

3. <https://python.sdv.univ-paris-diderot.fr/cours-python.pdf>

4. <http://sander.nabuurs.org/>

5. <http://www.cmbi.ru.nl/pythoncourse/spy/index.spy?site=python&action=Home>

6. <http://pgbovine.net/>

7. <http://pythontutor.com/>

1.3 Avant de commencer

Avant de débiter ce cours, voici quelques indications générales qui pourront vous servir pour la suite.

- Familiarisez-vous avec le site www.python.org⁸. Il contient énormément d'informations et de liens sur Python et vous permet en outre de le télécharger pour différentes plateformes (Linux, Mac, Windows). La page d'index des modules⁹ est particulièrement utile.
- Pour aller plus loin avec Python, Gérard Swinnen a écrit un très bon livre¹⁰ en français intitulé *Apprendre à programmer avec Python 3* et téléchargeable gratuitement. Les éditions Eyrolles proposent également la version papier¹¹ de cet ouvrage.
- Ce cours est basé sur une utilisation de Python sous Linux mais il est parfaitement transposable aux systèmes d'exploitation Windows et Mac.
- L'apprentissage d'un langage informatique comme Python va nécessiter d'écrire des lignes de codes à l'aide d'un éditeur de texte. Si vous êtes débutants, nous vous conseillons vivement d'utiliser *gedit* ou *geany*, qui sont les plus proches des éditeurs que l'on peut trouver sous Windows (*notepad*). Des éditeurs comme *emacs* et *vi* sont très puissants mais nécessitent un apprentissage particulier.

1.4 C'est quoi Python ?

Le langage de programmation Python a été créé en 1989 par Guido van Rossum, aux Pays-Bas. Le nom *Python* vient d'un hommage à la série télévisée *Monty Python's Flying Circus* dont G. van Rossum est fan. La première version publique de ce langage a été publiée en 1991.

La dernière version de Python est la version 3. Plus précisément, la version 3.7 a été publiée en juin 2018. Ce cours est basé sur Python 3. La version 2 de de Python est désormais obsolète, dans la mesure du possible évitez de l'utiliser.

La *Python Software Foundation* est l'association qui organise le développement de Python et anime la communauté de développeurs et d'utilisateurs.

Ce langage de programmation présente de nombreuses caractéristiques intéressantes : - Il est multiplateforme. C'est-à-dire qu'il fonctionne sur de nombreux systèmes d'exploitation : Windows, Mac OS, Linux, Android, iOS, depuis les mini-ordinateurs Raspberry Pi jusqu'aux supercalculateurs. - Il est gratuit. Vous pouvez l'installer sur autant d'ordinateur que vous voulez. - C'est un langage de haut niveau. Il demande relativement peu de connaissance sur le fonctionnement d'un ordinateur pour être utilisé. - C'est un langage interprété. Le programme n'a pas besoin de compiler son programme pour pouvoir l'utiliser, contrairement à des langages comme le C ou le C++. - Il est orienté objet. C'est-à-dire qu'il est possible de créer en Python des entités qui ont un sens dans le monde réel (une cellule, une protéine, un atome) avec un certains nombres de fonctionnement et d'interaction. - Il est relativement *simple* à prendre en main¹². - Enfin, il est très utilisé en bioinformatique et plus généralement en analyse de données.

Toutes ces caractéristiques font que Python est désormais enseigné dans de nombreuses formations, depuis l'enseignement secondaire jusqu'à l'enseignement supérieur.

8. <http://www.python.org>

9. <https://docs.python.org/3/py-modindex.html>

10. <http://www.inforef.be/swi/python.htm>

11. <http://www.eyrolles.com/Informatique/Livre/apprendre-a-programmer-avec-python-3-9782212134346>

12. Nous sommes d'accord, cette notion est très relative.

1.5 Premier contact avec Python sous Linux

Python est un langage interprété, c'est-à-dire que chaque ligne de code est lue puis interprétée afin d'être exécutée par l'ordinateur. Pour vous en rendre compte, lancez la commande :

```
1 python3
```

Celle-ci va lancer l'**interpréteur Python**. Vous devriez obtenir quelque chose de ce style :

```
1 pierre@jeera:~$ python3
2 Python 3.5.2 (default, Nov 17 2016, 17:05:23)
3 [GCC 5.4.0 20160609] on linux
4 Type "help", "copyright", "credits" or "license" for more information.
5 >>>
```

Le bloc `pierre@jeera:~$` représente l'invite de commande de votre *shell* sous Linux. Par la suite, cette invite de commande sera représentée simplement par le caractère `$`.

Le triple chevron `>>>` est l'invite de commande de Python (*prompt* en anglais), ce qui signifie que Python attend une commande. Tapez par exemple l'instruction

```
1 print("Hello world !")
```

puis validez votre commande en appuyant sur la touche **Entrée**.

Python a exécuté la commande directement et a affiché le texte `Hello world !`. Il attend ensuite votre prochaine instruction en affichant l'invite de l'interpréteur Python (`>>>`). En résumé, voici ce qui a du apparaître sur votre écran :

```
1 >>> print("Hello world !")
2 Hello world !
3 >>>
```

Vous pouvez refaire un nouvel essai en vous servant cette fois de l'interpréteur comme d'une calculatrice.

```
1 >>> 1+1
2 2
3 >>> 6*3
4 18
```

À ce stade, vous pouvez entrer une autre commande ou bien quitter l'interpréteur Python, soit en tapant la commande `exit()` puis en validant en appuyant sur la touche **Entrée**, soit en pressant simultanément les touches **Ctrl** et **D**.

Finalement l'interpréteur Python est un système interactif dans lequel vous pouvez entrer des commandes, que Python exécutera sous vos yeux (au moment où vous validerez la commande en tapant sur **Entrée**).

Il existe de nombreux autres langages interprétés tels que Perl¹³ ou R¹⁴. Le gros avantage est que l'on peut directement tester une commande à l'aide de l'interpréteur, ce qui est très utile pour déboguer (c'est-à-dire corriger les éventuelles erreurs d'un programme). Gardez bien en mémoire cette propriété de Python qui pourra parfois vous faire gagner un temps précieux !

13. <http://www.perl.org>

14. <http://www.r-project.org>

1.6 Premier programme Python

Bien sûr, l'interpréteur présente vite des limites dès lors que l'on veut exécuter une suite d'instructions plus complexe. Comme tout langage informatique, on peut enregistrer ces instructions dans un fichier, que l'on appelle communément un script Python.

Pour reprendre l'exemple précédent, ouvrez un éditeur de texte (par exemple *gedit* ou *nedit*) et entrez le code suivant.

```
1 print('Hello World !')
```

Ensuite enregistrez votre fichier sous le nom *test.py*, puis quittez l'éditeur de texte. L'extension standard des scripts Python est *.py*. Pour exécuter votre script, vous avez deux moyens.

1.6.1 Appel de l'interpréteur

Donnez le nom de votre script comme argument à la commande Python :

```
1 $ python3 test.py
2 Hello World !
3 $
```

1.6.2 Appel direct du script

Pour appeler directement votre script Python, deux opérations sont nécessaires :

1. Précisez au *shell* la localisation de l'interpréteur Python en indiquant dans la première ligne du script :

```
1 #! /usr/bin/env python3
```

Dans notre exemple, le script *test.py* contient alors le texte suivant :

```
1 #!/usr/bin/env python3
2
3 print('Hello World !')
```

2. Rendez votre script Python exécutable en tapant :

```
1 chmod +x test.py
```

Pour exécuter votre script, tapez son nom précédé des deux caractères *./* (afin de préciser au *shell* où se trouve votre script) :

```
1 $ ./test.py
2 Hello World !
3 $
```

1.7 Commentaires

Dans un script, tout ce qui suit le caractère `#` est ignoré par Python jusqu'à la fin de la ligne et est considéré comme un commentaire. Une exception notable est la première ligne de votre script qui peut être

```
1 #! /usr/bin/env python3
```

et qui a alors une signification particulière pour Python.

Les commentaires sont indispensables pour que vous puissiez annoter votre code. Il faut absolument les utiliser pour décrire les principales parties de votre code dans un langage humain.

```
1 #! /usr/bin/env python3
2
3 # votre premier script Python
4 print('Hello World !')
5
6 # d'autres commandes plus utiles pourraient suivre
```

1.8 Notion de bloc d'instructions et d'indentation

Pour terminer ce chapitre d'introduction, nous allons aborder dès maintenant les notions de **bloc d'instructions** et d'**indentation**.

En programmation, il est courant de répéter un certain nombre de choses (avec les boucles, voir chapitre 5) ou de faire des choix (avec les tests, voir chapitre 6).

Par exemple, imaginez que vous souhaitiez répéter 10 fois 3 instructions différentes, les unes à la suite des autres. Regardez l'exemple suivant en pseudo-code :

```
1 instruction_qui_indique_à_Python_de_répéter_10_fois_ce_qui_suit:
2     sous_instruction1
3     sous_instruction2
4     sous_instruction3
5 instruction_suivante
```

La première ligne correspond à une instruction qui va indiquer à Python de répéter 10 fois d'autres instructions (il s'agit d'une boucle, on verra le nom de la commande exacte plus tard). Dans le pseudo-code ci-dessus, il y a 3 instructions à répéter, nommées `sous_instruction1` à `sous_instruction3`.

Notez bien les détails de la syntaxe. La première ligne indique que l'on veut répéter une ou plusieurs instructions, elle se termine par `:`. Ce symbole `:` indique à Python qu'il doit attendre un bloc d'instructions, c'est-à-dire un certains nombres de sous-instructions à répéter. Python reconnaît un bloc d'instructions car il est indenté. L'indentation est le décalage d'un ou plusieurs espaces ou tabulations des instructions `sous_instruction1` à `sous_instruction3`, par rapport à `instruction_qui_indique_à_Python_de_répéter_10_fois_ce_qui_suit:`.

Pratiquement, l'indentation en Python doit être homogène (soit des espaces, soit des tabulations, mais pas un mélange des deux). Une indentation avec 4 espaces est le style d'indentation le plus traditionnel et celui qui permet une bonne lisibilité du code.

Enfin, la ligne `instruction_suivante` sera exécutée une fois que le bloc d'instructions sera terminé.

Si tout cela vous semble un peu fastidieux, ne vous inquiétez pas. Vous allez comprendre tous ces détails, et surtout les acquérir, en continuant ce cours chapitre par chapitre.

1.9 Python 2 ou Python 3 ?

Ce cours est basé sur la version 3 de Python, qui est maintenant devenu un standard. La plupart des projets importants ont migré vers Python 3 ¹⁵. Par ailleurs Python 2 cessera d'être maintenu au delà de mi-2020 ¹⁶. Python 3 est donc la version que nous vous recommandons fortement.

Si néanmoins, vous devez un jour travailler sur un ancien programme écrit en Python 2, sachez qu'il existe quelques différences importantes entre Python 2 et Python 3. Le chapitre 19 *Pour aller plus loin* vous donnera plus de précisions.

15. <http://www.python3statement.org>

16. <https://pythonclock.org/>

Chapitre 2

Variables

2.1 Définition d'une variable

Une **variable** est une zone de la mémoire dans laquelle une **valeur** est stockée. Aux yeux du programmeur, cette variable est définie par un **nom**, alors que pour l'ordinateur, il s'agit en fait d'une adresse (*i.e.* une zone particulière de la mémoire).

En Python, la **déclaration** d'une variable et son **initialisation** (c'est-à-dire la première valeur que l'on va stocker dedans) se font en même temps. Pour vous en convaincre, testez les instructions suivantes après avoir lancé l'interpréteur :

```
1 >>> x = 2
2 >>> x
3 2
```

Dans cet exemple, nous avons déclaré, puis initialisé la variable `x` avec la valeur `2`. Notez bien qu'en réalité, il s'est passé plusieurs choses :

- Python a *deviné* que la variable était un entier. On dit que Python est un langage au *typage dynamique*.
- Python a alloué (*i.e.* réservé) l'espace en mémoire pour y accueillir un entier (chaque type de variable prend plus ou moins d'espace en mémoire), et a fait en sorte qu'on puisse retrouver la variable sous le nom `x`
- Python a assigné la valeur `2` à la variable `x`.

Dans certains autres langages, il faut coder ces différentes étapes une par une (en C par exemple). Python étant un langage dit de *haut niveau*, la simple instruction `x = 2` a suffi à réaliser les 3 étapes en une fois !

Ensuite, l'interpréteur nous a permis de connaître le contenu de la variable juste en tapant son nom. Retenez ceci car c'est une **spécificité de l'interpréteur Python**, très pratique pour chasser les erreurs (*debugging*) dans un programme. Par contre, la ligne d'un script Python qui contient seulement le nom d'une variable (sans aucune autre indication) n'affichera pas la valeur de la variable à l'écran (pour autant ceci reste valide et ne générera pas d'erreur).

Dernière chose, l'opérateur d'affectation `=` s'utilise dans un certain sens : par exemple `x = 2` signifie qu'on attribue la valeur située à droite de l'opérateur `=` (`2`) à la variable située à gauche (`x`). Certains autres langages comme **R** utilise les symboles `<-` pour rendre les choses plus explicites, par exemple `x <- 2`.

Si on a `x = y - 3`, l'opération `y - 3` est d'abord évaluée et ensuite le résultat de cette opération est affecté à la variable `x`.

2.2 Les types de variables

Le **type** d'une variable correspond à la nature de celle-ci. Les trois types principaux dont nous aurons besoin dans un premier temps sont les entiers (*integer* ou *int*), les réels (*float*) et les chaînes de caractères (*string* ou *str*). Bien sûr, il existe de nombreux autres types (par exemple, les nombres complexes), c'est d'ailleurs un des gros avantages de Python (si vous n'êtes pas effrayés, vous pouvez vous en rendre compte ici¹).

Dans l'exemple précédent, nous avons stocké un nombre entier (*int*) dans la variable `x`, mais il est tout à fait possible de stocker des nombres réels (*float*), des chaînes de caractères (*string* ou *str*) ou plein d'autres types de variables que nous verrons par la suite :

```
1 >>> y = 3.14
2 >>> y
3 3.14
4 >>> a = "bonjour"
5 >>> a
6 'bonjour'
7 >>> b = 'salut'
8 >>> b
9 'salut'
10 >>> c = '''girafe'''
11 >>> c
12 'girafe'
```

Vous remarquez que Python reconnaît certains types de variable automatiquement (entier, réel). Par contre, pour une chaîne de caractères, il faut l'entourer de guillemets (simples, doubles voire trois guillemets successifs simples ou doubles) afin d'indiquer à Python le début et la fin de la chaîne.

2.3 Nommage des variables

Le nom des variable en Python peut-être constitué de lettres minuscules (a à z), de lettres majuscules (A à Z), de nombres (0 à 9) ou du caractère souligné (_).

Néanmoins, un nom de variable ne doit pas débiter ni par un chiffre, ni par `_` et ne peut pas contenir de caractère accentué. Il faut absolument éviter d'utiliser un mot "réservé" par Python comme nom de variable (par exemple : `print`, `range`, `for`, `from`, etc.).

Python est sensible à la casse, ce qui signifie que les variables `Test`, `test` ou `TEST` sont différentes. Enfin, vous ne pouvez pas utiliser d'espace dans un nom de variable.

2.4 Opérations

2.4.1 Opérations sur les types numériques

Les quatre opérations de base se font de manière simple sur les types numériques (nombres entiers et réels) :

```
1 >>> x = 45
2 >>> x + 2
3 47
```

1. <https://docs.python.org/3.6/library/stdtypes.html?highlight=type#>

```

4 >>> y = 2.5
5 >>> x + y
6 47.5
7 >>> (x * 10) / y
8 180.0

```

Remarquez toutefois que si vous mélangez les types entiers et réels, le résultat est renvoyé comme un réel (car ce type est plus général). Par ailleurs, l'utilisation de parenthèses permet de gérer les priorités.

L'opérateur puissance utilise le symbole **. Pour obtenir le reste d'une division entière (voir ici² pour un petit rappel sur la division entière), on utilise le symbole modulo % :

```

1 >>> 2**3
2 8
3 >>> 5 % 4
4 1
5 >>> 8 % 4
6 0

```

Les symboles +, -, *, /, ** et % sont appelés **opérateurs**, car ils permettent de faire des opérations sur les variables.

2.4.2 Opérations sur les chaînes de caractères

Pour les chaînes de caractères, deux opérations sont possibles, l'addition et la multiplication :

```

1 >>> chaine = "Salut"
2 >>> chaine
3 'Salut'
4 >>> chaine + " Python"
5 'Salut Python'
6 >>> chaine * 3
7 'SalutSalutSalut'

```

L'opérateur d'addition + permet de concaténer (assembler) deux chaînes de caractères et l'opérateur de multiplication * permet de dupliquer plusieurs fois une chaîne.

Attention : Vous voyez que les opérateurs + et * se comportent différemment selon s'il s'agit d'entiers ou de chaînes de caractères : 2 + 2 est une addition, '2' + '2' est une concaténation. On appelle ce comportement **surcharge des opérateurs**. Nous serons amenés à revoir cette notion dans le chapitre 18 sur les classes.

2.4.3 Opérations illicites

Attention à ne pas faire d'opération illicite car vous obtiendriez un message d'erreur :

```

1 >>> 'toto' + 2
2 Traceback (most recent call last):
3   File "<stdin>", line 1, in ?
4 TypeError: cannot concatenate 'str' and 'int' objects

```

2. https://fr.wikipedia.org/wiki/Division_euclidienne

Notez que Python vous donne le maximum d'information dans son message d'erreur. Dans l'exemple précédent, il vous indique que vous ne pouvez pas mélanger des objets de type `str` (*string*, donc des chaînes de caractères) avec des objets de type `int` (donc des entiers), ce qui est assez logique.

2.5 La fonction `type()`

Si vous ne vous souvenez plus du type d'une variable, utilisez la fonction `type()` qui vous le rappellera.

```
1 >>> x = 2
2 >>> type(x)
3 <class 'int'>
4 >>> y = 2.0
5 >>> type(y)
6 <class 'float'>
7 >>> z = '2'
8 >>> type(z)
9 <class 'str'>
```

Faites bien attention, car pour Python, la valeur `2` (nombre entier) est différente de `2.0` (nombre réel), de même que `2` (nombre entier) est différent de `'2'` (chaîne de caractères). Nous verrons plus tard ce que signifie le mot *class*.

2.6 Conversion de types

Dans tout langage de programmation, on est souvent amené à convertir les types, c'est-à-dire passer d'un type numérique à une chaîne de caractères ou vice-versa. En Python, rien de plus simple avec les fonctions `int()`, `float()` et `str()`. Pour vous en convaincre, regardez ces exemples :

```
1 >>> i = 3
2 >>> str(i)
3 '3'
4 >>> i = '456'
5 >>> int(i)
6 456
7 >>> float(i)
8 456.0
9 >>> i = '3.1416'
10 >>> float(i)
11 3.1416
```

On verra au chapitre 7 sur les fichiers que ces conversions sont essentielles. En effet, lorsqu'on lit ou écrit des nombres dans un fichier, ils sont considérés comme du texte.

Toute conversion d'une variable d'un type en un autre est appelé *casting* en anglais, il se peut que vous croisiez ce terme si vous allez consulter d'autres ressources.

2.7 Note sur la division

Notez bien qu'en Python 3, la division de nombres entiers renvoie par défaut un nombre réel (*float*) :

```
1 >>> x = 3 / 4
2 >>> x
3 0.75
4 >>> type(x)
5 <class 'float'>
```

Attention, ceci n'était pas le cas en Python 2. Pour en savoir plus sur ce point, vous pouvez consulter la section *Pour aller plus loin*

2.8 Note sur le vocabulaire et la syntaxe

Nous avons vu dans ce chapitre la notion de **variable** qui est commune à tous les langages de programmation. Toutefois, Python est un langage dit **orienté objet**, il se peut que dans la suite du cours nous employions ce mot *objet* pour désigner une variable. Par exemple *variable de type entier* sera équivalent à un *objet de type entier*. Nous verrons ce que le mot *objet* signifie réellement plus tard (tout comme le mot *class*).

Par ailleurs, nous avons rencontré plusieurs fois des **fonctions** dans ce chapitre, avec `type(x)`, `int(x)`, `float(x)` et `str(x)`. Dans le chapitre 1 nous avons également vu la fonction `print()`. On reconnaît qu'il s'agit d'une fonction au nom - par exemple `type` - suivi de parenthèses (). En Python la syntaxe générale est `fonction()`. La variable `x` entre les parenthèses est appelé **argument** que l'on passe à la fonction. Dans `type(2)` c'est l'entier 2 qui est l'argument passé à la fonction. Pour l'instant on retiendra qu'une fonction est une sorte de *boîte* à qui on passe un *argument* et qui renvoie un *résultat* ou plus généralement un objet. Par exemple, la fonction `type()` renvoie le type de la variable qu'on lui a passé en argument.

Si ces notions vous font peur, ne vous inquiétez pas, au fur et à mesure que vous avancerez dans le cours tout deviendra limpide.

Chapitre 3

Affichage

Nous avons déjà vu au chapitre 1 la fonction `print()` qui permet d'afficher une chaîne de caractères. Elle permet en plus d'afficher le contenu d'une ou plusieurs variables :

```
1 >>> x = 32
2 >>> nom = 'John'
3 >>> print(nom , ' a ' , x , ' ans')
4 John a 32 ans
```

Python a donc écrit la phrase en remplaçant les variables `x` et `nom` par leur contenu. Vous pouvez noter également que pour écrire plusieurs blocs de texte sur une seule ligne, nous avons utilisé le séparateur `,` avec la fonction `print()`. En regardant de plus près, vous vous apercevrez que Python a automatiquement ajouté un espace à chaque fois que l'on utilisait le séparateur `,`. Par conséquent, si vous voulez mettre un seul espace entre chaque bloc, vous pouvez retirer ceux de vos chaînes de caractères :

```
1 >>> print(nom , 'a' , x , 'ans')
2 John a 32 ans
```

Pour imprimer deux chaînes de caractères l'une à côté de l'autre sans espace, vous devrez les concaténer :

```
1 >>> ani1 = 'chat'
2 >>> ani2 = 'souris'
3 >>> print(ani1, ani2)
4 chat souris
5 >>> print(ani1 + ani2)
6 chatsouris
```

3.1 Écriture formatée

La méthode `.format()` permet une meilleure organisation de l'affichage des variables (nous expliquerons à la fin de ce chapitre ce que signifie le terme *méthode* en Python).

Si on reprend l'exemple précédent :

```
1 >>> x = 32
2 >>> nom = 'John'
3 >>> print('{} a {} ans'.format(nom, x))
4 John a 32 ans
```

- Dans la chaîne de caractères, les accolades vides {} précisent l'endroit où le contenu de la variable doit être inséré.
- Juste après la chaîne de caractères, l'instruction `.format(nom, x)` indique la liste des variables à insérer, d'abord la variable `nom` puis la variable `x`. Ainsi, la méthode `.format()` agit sur la chaîne de caractères à laquelle elle est attachée par le ..

Remarque : il est possible d'indiquer entre les accolades {} dans quel ordre afficher les variables, avec 0 pour la variable à afficher en premier, 1 pour la variable à afficher en second, etc. (attention, Python commence à compter à 0). Cela permet de modifier l'ordre dans lequel sont affichées les variables.

```

1 >>> x = 32
2 >>> nom = 'John'
3 >>> print('{0} a {1} ans'.format(nom, x))
4 John a 32 ans
5 >>> print('{1} a {0} ans'.format(nom, x))
6 32 a John ans

```

Imaginez maintenant que vous vouliez calculer puis afficher la proportion de GC d'un génome. La proportion de GC s'obtient comme la somme des bases Guanine (G) et Cytosine (C) divisée par le nombre total de bases (A, T, C, G) du génome considéré. Si on a, par exemple, 4500 bases G, 2575 bases C pour un total de 14800 bases, vous pourriez faire comme suit (notez bien l'utilisation des parenthèses pour gérer les priorités des opérateurs) :

```

1 >>> propGC = (4500 + 2575) / 14800
2 >>> print("La proportion de GC est", propGC)
3 La proportion de GC est 0.4780405405405405

```

Le résultat obtenu présente trop de décimales (seize dans le cas présent)... Pour écrire le résultat plus lisiblement, vous pouvez spécifier dans les accolades {} le format qui vous intéresse. Dans votre cas, vous voulez formater un réel (*float*) pour l'afficher avec deux puis trois décimales :

```

1 >>> print("La proportion de GC est {:.2f}".format(propGC))
2 Le proportion de GC est 0.48
3 >>> print("La proportion de GC est {:.3f}".format(propGC))
4 La proportion de GC est 0.478

```

Détaillons le contenu des accolades :

- Les deux points : indiquent que l'on veut préciser le format.
- La lettre *f* indique que l'on souhaite afficher la variable sous forme d'un réel (*float*).
- Les caractères *.2* indiquent la précision voulue, soit ici deux chiffres après la virgule.
- Notez enfin que le formatage avec *.xf* (*x* étant un entier positif) renvoie un résultat arrondi.

Il est par ailleurs possible de combiner le formatage (à droite des 2 points) ainsi que l'emplacement des variables à substituer (à gauche des 2 points), par exemple :

```

1 >>> print("propGC(2 décimales) = {0:.2f}, propGC(3 décimales) = {0:.3f}".format(propGC))
2 propGC(2 décimales) = 0.48, propGC(3 décimales) = 0.478

```

Vous remarquerez qu'on utilise ici la même variable (`propGC`) à deux endroits différents.

Vous pouvez aussi formater des entiers avec *d*,

```

1 >>> nbG = 4500
2 >>> print("Le génome de cet exemple contient {:d} guanines".format(nbG))
3 Le génome de cet exemple contient 4500 guanines

```


ou mettre plusieurs nombres dans une même chaîne de caractères.

```

1 >>> nbG = 4500
2 >>> nbC = 2575
3 >>> print("Ce génome contient {:d} G et {:d} C, soit une prop de GC de {:.2f}" \
4 ... .format(nbG,nbC,propGC))
5 Ce génome contient 4500 G et 2575 C, soit une prop de GC de 0.48
6 >>> percGC = propGC * 100
7 >>> print "Ce génome contient {:d} G et {:d} C, soit un %GC de {:.2f} %" \
8 ... .format(nbG,nbC,percGC)
9 Ce génome contient 4500 G et 2575 C, soit un %GC de 47.80 %

```

Remarque : Le signe `\` en fin de ligne permet de poursuivre la commande sur la ligne suivante. Cette syntaxe est pratique lorsque vous voulez taper une commande longue.

Enfin, il est possible de préciser sur combien de caractères vous voulez qu'un résultat soit écrit et comment se fait l'alignement (à gauche, à droite ou centré). Dans la portion de code suivant, le caractère `;` sert de séparateur entre les instructions sur une même ligne :

```

1 >>> print(10) ; print(1000)
2 10
3 1000
4 >>> print("{:>6d}".format(10)) ; print("{:>6d}".format(1000))
5     10
6     1000
7 >>> print("{:<6d}".format(10)) ; print("{:<6d}".format(1000))
8 10
9 1000
10 >>> print("{:^6d}".format(10)) ; print("{:^6d}".format(1000))
11  10
12  1000
13 >>> print("{:*^6d}".format(10)) ; print("{:*^6d}".format(1000))
14 **10**
15 *1000*
16 >>> print("{:0>6d}".format(10)) ; print("{:0>6d}".format(1000))
17 000010
18 001000

```

Notez que `>` spécifie un alignement à droite, `<` spécifie un alignement à gauche et `^` spécifie un alignement centré. Il est également possible d'indiquer le caractère qui servira de remplissage lors des alignements (l'espace par défaut).

Ce formatage est également possible sur des chaînes de caractères, notées `s` (comme *string*) :

```

1 >>> print "atom HN" ; print "atom HDE1"
2 atom HN
3 atom HDE1
4 >>> print "atom {:>4s}".format("HN") ; print "atom {:>4s}".format("HDE1")
5 atom  HN
6 atom  HDE1

```

Vous voyez tout de suite l'énorme avantage de l'écriture formatée. Cela vous permet d'écrire en colonnes parfaitement alignées. Nous verrons que ceci est très pratique si l'on veut écrire les coordonnées des atomes d'une molécule au format PDB.

Pour les réels, il est possible de combiner le nombre de chiffres après la virgule.

```

1 >>> print("{:7.3f}".format(propGC))
2 47.804
3 >>> print("{:10.3f}".format(propGC))
4 47.804

```

L'instruction `7.3f` signifie que l'on souhaite écrire le réel avec 3 décimales et formaté sur 7 caractères (par défaut justifiés à droite). L'instruction `10.3f` fait la même chose sur 10 caractères. Remarquez que le séparateur décimal `.` compte pour un caractère.

Enfin, si on veut écrire des accolades littérales et utiliser la méthode `format()` en même temps, il faudra doubler les accolades pour échapper au formatage.

```

1 >>> print("Accolades littérales {} et accolades pour le formatage {}".format(10))
2 Accolades littérales {} et accolades pour le formatage 10

```

À retenir : Comme indiqué ci-dessus, la méthode `format()` agit sur la chaîne de caractères à laquelle elle est *attachée* par un `.`, et n'a rien à voir avec la fonction `print()`. Si on donne une chaîne suivie d'un `.format()` à la fonction `print()`, Python évalue d'abord le formatage et c'est la chaîne qui en résulte qui est affichée à l'écran. Tout comme dans `print(5*5)`, c'est d'abord la multiplication qui est évaluée puis son résultat qui est affiché à l'écran. On peut s'en rendre compte de la manière suivante dans l'interpréteur :

```

1 >>> "{:10.3f}".format(propGC)
2 ' 47.804'
3 >>> type("{:10.3f}".format(propGC))
4 <class 'str'>

```

Python affiche le résultat du `format()` comme une chaîne de caractères (*string*), et la fonction `type()` nous le confirme.

3.2 Ancienne méthode de formatage des chaînes de caractères

Conseil : Pour les débutants, vous pouvez passer cette section.

Sur d'anciens programmes Python, il se peut que vous rencontriez l'écriture formatée dans le style suivant :

```

1 >>> x = 32
2 >>> nom = 'John'
3 >>> print("%s a %d ans" % (nom,x))
4
5 >>> nbG = 4500
6 >>> nbC = 2575
7 >>> propGC = (4500.0 + 2575)/14800
8 >>> print("On a %d G et %d C -> prop GC = %.2f" % (nbG,nbC,propGC))
9 On a 4500 G et 2575 C -> prop GC = 0.48

```

La syntaxe est légèrement différente. Le symbole `%` est appelé une première fois dans la chaîne de caractères (dans l'exemple ci-dessus `%d`, `%d` et `%.2f`) pour :

- Désigner l'endroit où sera placé la variable dans la chaîne de caractères.
- Préciser le type de variable à formater, `d` pour un entier (i fonctionne également) ou `f` (*float*) pour un réel.

- Eventuellement pour indiquer le format voulu. Ici `.2` signifie une précision de deux chiffres après la virgule.

Le signe `%` est rappelé une seconde fois (`% (nbG,nbC,propGC)`) pour indiquer les variables à formater. Cette ancienne façon de formater vous est présentée à titre d'information. Ne l'utilisez pas dans vos programmes.

3.3 Note sur le vocabulaire et la syntaxe

Revenons quelques instants sur la notion de **méthode** abordé dans ce chapitre avec `format()`. En Python, on peut finalement considérer chaque variable comme un objet sur lequel on peut appliquer des méthodes. Une méthode est simplement une fonction qui utilise et/ou agit sur l'objet lui-même, les deux étant connectés par un point `.`. La syntaxe générale est du type `objet.méthode()`.

Dans l'exemple suivant :

```
1 >>> "Joe a {} ans".format(20)
2 'Joe a 20 ans'
```

la méthode `format()` est lié à `"Joe a {} ans"` qui est un objet de type chaîne de caractères (*string*). La méthode renvoie une nouvelle chaîne de caractères avec le bon formatage.

Nous aurons de nombreuses occasions de revoir cette notation `objet.méthode()`.

3.4 Exercices

Conseil : utilisez l'interpréteur Python pour les exercices 2 à 4.

3.4.1 Affichage dans l'interpréteur et dans un programme

Ouvrez l'interpréteur Python et tapez `1+1`. Que se passe-t-il? Ecrivez la même chose dans un script `test.py`. Exécutez ce script en tapant `python3 test.py` dans un *shell* Unix. Que se passe-t-il? Pourquoi? Faites en sorte d'afficher le résultat de l'addition `1+1` en exécutant le script dans un *shell* Unix.

3.4.2 Poly-A

Générez une chaîne de caractères représentant un oligonucléotide polyA (AAAA...) de 20 bases de longueurs, sans taper littéralement toutes les bases.

3.4.3 Poly-A et poly-GC

Suivant le modèle du dessus, générez en une ligne de code un polyA de 20 bases suivi d'un polyGC régulier (GCGCGC...) de 40 bases.

3.4.4 Écriture formatée

En utilisant l'écriture formatée, affichez en une seule ligne les variables `a`, `b` et `c` dont les valeurs sont respectivement `"salut"`, `102` et `10.318`.

3.4.5 Écriture formatée 2

Dans un script `propGC.py`, calculez un pourcentage de GC avec l'instruction suivante : `percGC = ((4500 + 2575)/14800)*100`. Ensuite, affichez le contenu de la variable `percGC` à l'écran avec 0, 1, 2, 3 puis 4 décimales sous forme arrondie en utilisant `format()`. On souhaite que le programme affiche la sortie suivante :

```
1 Le pourcentage de GC est 48 %
2 Le pourcentage de GC est 47.8 %
3 Le pourcentage de GC est 47.80 %
4 Le pourcentage de GC est 47.804 %
5 Le pourcentage de GC est 47.8041 %
```

Chapitre 4

Listes

4.1 Définition

Une **liste** est une structure de données qui contient une série de valeurs. Python autorise la construction de liste contenant des valeurs de type différent (par exemple entier et chaîne de caractères), ce qui leur confère une grande flexibilité. Une liste est déclarée par une série de valeurs (ne pas oublier les guillemets, simples ou doubles, s'il s'agit de chaînes de caractères) séparées par des **virgules**, et le tout encadré par des **crochets**. En voici quelques exemples :

```
1 >>> animaux = ['girafe','tigre','singe','souris']
2 >>> tailles = [5, 2.5, 1.75, 0.15]
3 >>> mixte = ['girafe', 5, 'souris', 0.15]
4 >>> animaux
5 ['girafe', 'tigre', 'singe', 'souris']
6 >>> tailles
7 [5, 2.5, 1.75, 0.15]
8 >>> mixte
9 ['girafe', 5, 'souris', 0.15]
```

Lorsque l'on affiche une liste, Python la restitue telle qu'elle a été saisie.

4.2 Utilisation

Un des gros avantages d'une liste est que vous pouvez appeler ses éléments par leur position. Ce numéro est appelé **indice** (ou index) de la liste.

```
1 liste : ['girafe', 'tigre', 'singe', 'souris']
2 indice :      0      1      2      3
```

Soyez très **attentifs** au fait que les indices d'une liste de n éléments commence à 0 et se termine à $n-1$. Voyez l'exemple suivant :

```
1 >>> animaux = ['girafe','tigre','singe','souris']
2 >>> animaux[0]
3 'girafe'
4 >>> animaux[1]
5 'tigre'
6 >>> animaux[3]
7 'souris'
```

Par conséquent, si on appelle l'élément d'indice 4 de notre liste, Python renverra un message d'erreur :

```

1 >>> animaux[4]
2 Traceback (innermost last):
3   File "<stdin>", line 1, in ?
4 IndexError: list index out of range

```

N'oubliez pas ceci ou vous risqueriez d'obtenir des bugs inattendus!

4.3 Opération sur les listes

Tout comme les chaînes de caractères, les listes supportent l'opérateur + de concaténation, ainsi que l'opérateur * pour la duplication :

```

1 >>> ani1 = ['girafe','tigre']
2 >>> ani2 = ['singe','souris']
3 >>> ani1 + ani2
4 ['girafe', 'tigre', 'singe', 'souris']
5 >>> ani1 * 3
6 ['girafe', 'tigre', 'girafe', 'tigre', 'girafe', 'tigre']

```

L'opérateur + est très pratique pour concaténer deux listes. Vous pouvez aussi utiliser la fonction `append()` lorsque vous souhaitez ajouter un seul élément à la fin d'une liste.

Dans l'exemple suivant nous allons créer une liste vide puis lui ajouter deux éléments, d'abord avec la concaténation :

```

1 >>> a = [] # liste vide
2 >>> a
3 []
4 >>> a = a + [15]
5 >>> a
6 [15]
7 >>> a = a + [-5]
8 >>> a
9 [15, -5]

```

puis avec la fonction `append()` :

```

1 >>> b = []
2 >>> b
3 []
4 >>> b.append(15)
5 >>> b
6 [15]
7 >>> b.append(-5)
8 >>> b
9 [15, -5]

```

Dans l'exemple ci-dessus nous obtenons le même résultat pour faire grandir une liste en utilisant l'opérateur de concaténation ou la fonction `.append()`. Nous vous conseillons néanmoins dans ce cas précis d'utiliser la fonction `.append()` dont la syntaxe est plus élégante.

Nous reverrons en détail la fonction `.append()` dans le chapitre 11 - *Plus sur les listes*, ainsi que la notation `objet.méthode()` dans les chapitres suivants.

4.4 Indichage négatif

La liste peut également être indexée avec des nombres négatifs selon le modèle suivant :

```

1 liste      : ['girafe', 'tigre', 'singe', 'souris']
2 indice positif :      0      1      2      3
3 indice négatif :     -4     -3     -2     -1

```

ou encore :

```

1 liste      : ['A','B','C','D','E','F']
2 indice positif : 0  1  2  3  4  5
3 indice négatif : -6 -5 -4 -3 -2 -1

```

Les indices négatifs reviennent à compter à partir de la fin. Leur principal avantage est que vous pouvez accéder au dernier élément d'une liste à l'aide de l'indice `-1` sans pour autant connaître la longueur de cette liste. L'avant-dernier élément a lui l'indice `-2`.

```

1 >>> animaux = ['girafe','tigre','singe','souris']
2 >>> animaux[-1]
3 'souris'
4 >>> animaux[-2]
5 'singe'

```

Pour accéder au premier élément de la liste, il faut par contre connaître le bon indice :

```

1 >>> animaux[-4]
2 'girafe'

```

Dans ce cas, on utiliserait plutôt `animaux[0]`.

4.5 Tranches

Un autre avantage des listes est la possibilité de sélectionner une partie en utilisant un indichage construit sur le modèle `[m:n+1]` pour récupérer tous les éléments, du *m*ème au *n*ème (de l'élément *m* inclus à l'élément *n+1* exclus). On dit alors qu'on récupère une **tranche** de la liste, par exemple :

```

1 >>> animaux = ['girafe', 'tigre', 'singe', 'souris']
2 >>> animaux[0:2]
3 ['girafe', 'tigre']
4 >>> animaux[0:3]
5 ['girafe', 'tigre', 'singe']
6 >>> animaux[0:]
7 ['girafe', 'tigre', 'singe', 'souris']
8 >>> animaux[:]
9 ['girafe', 'tigre', 'singe', 'souris']
10 >>> animaux[1:]
11 ['tigre', 'singe', 'souris']
12 >>> animaux[1:-1]
13 ['tigre', 'singe']

```

Notez que lorsqu'aucun indice n'est indiqué à gauche ou à droite du symbole `:`, Python prend par défaut tous les éléments depuis le début ou tous les éléments jusqu'à la fin respectivement.

On peut aussi préciser le pas en ajoutant un `:` supplémentaire et en indiquant le pas par un entier.

```
1 >>> animaux = ['girafe', 'tigre', 'singe', 'souris']
2 >>> animaux[0:3:2]
3 ['girafe', 'singe']
4 >>> x = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
5 >>> x
6 [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
7 >>> x[::1]
8 [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
9 >>> x[::2]
10 [0, 2, 4, 6, 8]
11 >>> x[::3]
12 [0, 3, 6, 9]
13 >>> x[1:6:3]
14 [1, 4]
```

Finalement, on voit que l'accès au contenu d'une liste avec des crochets fonctionne sur le modèle `liste[début:fin:pas]`.

4.6 Fonction len()

L'instruction `len()` vous permet de connaître la longueur d'une liste, c'est-à-dire le nombre d'éléments que contient la liste. Voici un exemple d'utilisation :

```
1 >>> animaux = ['girafe', 'tigre', 'singe', 'souris']
2 >>> len(animaux)
3 4
4 >>> len([1, 2, 3, 4, 5, 6, 7, 8])
5 8
```

4.7 Les fonctions range() et list()

L'instruction `range()` est une fonction spéciale en Python qui va nous permettre de générer des nombres entiers compris dans un intervalle lorsqu'elle est utilisée en combinaison avec la fonction `list()`. Par exemple :

```
1 >>> list(range(10))
2 [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
```

La commande `list(range(10))` a généré une liste contenant tous les nombres entiers de 0 inclus à 10 **exclus**. Nous verrons l'utilisation de la fonction `range()` toute seule dans le chapitre suivant sur les boucles.

Dans l'exemple ci-dessus, la fonction `range()` a pris un argument, mais elle peut également prendre deux ou trois arguments, voyez plutôt :

```
1 >>> list(range(0,5))
2 [0, 1, 2, 3, 4]
3 >>> list(range(15,20))
4 [15, 16, 17, 18, 19]
5 >>> list(range(0,1000,200))
```



```

6 [0, 200, 400, 600, 800]
7 >>> list(range(2,-2,-1))
8 [2, 1, 0, -1]

```

L'instruction `range()` fonctionne sur le modèle `range([début,] fin[, pas])`. Les arguments entre crochets sont optionnels. Pour obtenir une liste, il faut l'utiliser systématiquement avec la fonction `list()`.

Enfin, prenez garde aux arguments optionnels par défaut (0 pour début et 1 pour pas) :

```

1 >>> list(range(10,0))
2 []

```

Ici la liste est vide car Python a pris la valeur du pas par défaut qui est de 1. Ainsi, si on commence à 10 et qu'on avance par pas de 1, on ne pourra jamais atteindre 0. Python génère ainsi une liste vide. Pour éviter ça, il faudra absolument préciser le pas de -1 pour les listes décroissantes :

```

1 >>> list(range(10,0,-1))
2 [10, 9, 8, 7, 6, 5, 4, 3, 2, 1]

```

4.8 Listes de listes

Pour finir, sachez qu'il est tout-à-fait possible de construire des listes de listes. Cette fonctionnalité peut être parfois très pratique. Par exemple :

```

1 >>> enclos1 = ['girafe', 4]
2 >>> enclos2 = ['tigre', 2]
3 >>> enclos3 = ['singe', 5]
4 >>> zoo = [enclos1, enclos2, enclos3]
5 >>> zoo
6 [['girafe', 4], ['tigre', 2], ['singe', 5]]

```

Dans cet exemple, chaque sous-liste contient une catégorie d'animal et le nombre d'animaux pour chaque catégorie.

Pour accéder à un élément de la liste, on utilise l'indiciage habituel :

```

1 >>> zoo[1]
2 ['tigre', 2]

```

Pour accéder à un élément de la sous-liste, on utilise un double indiciage :

```

1 >>> zoo[1][0]
2 'tigre'
3 >>> zoo[1][1]
4 2

```

On verra un peu plus loin qu'il existe en Python les dictionnaires qui sont très pratiques pour faire ce genre de choses. On verra aussi qu'il existe un module nommé `numpy` permettant de gérer des listes ou tableaux de nombres (vecteurs et matrices), ainsi que de faire des opérations avec.

4.9 Exercices

Conseil : utilisez l'interpréteur Python.

4.9.1 Jours de la semaine

Constituez une liste `semaine` contenant les 7 jours de la semaine.

1. À partir de cette liste, comment récupérez-vous seulement les 5 premiers jours de la semaine d'une part, et ceux du week-end d'autre part (*utilisez pour cela l'indiciage*) ?
2. Cherchez un autre moyen pour arriver au même résultat (*en utilisant un autre indiciage*).
3. Trouvez deux manières pour accéder au dernier jour de la semaine.
4. Inversez les jours de la semaine en une commande.

4.9.2 Saisons

Créez 4 listes `hiver`, `printemps`, `ete` et `automne` contenant les mois correspondants à ces saisons. Créez ensuite une liste `saisons` contenant les sous-listes `hiver`, `printemps`, `ete` et `automne`. Prévoyez ce que valent les variables suivantes, puis vérifiez-le dans l'interpréteur :

1. `saisons[2]`
2. `saisons[1][0]`
3. `saisons[1:2]`
4. `saisons[:][1]`. Comment expliquez-vous ce dernier résultat ?

4.9.3 Table des 9

Affichez la table des 9 en une seule commande avec les instructions `range()` et `list()`.

4.9.4 Nombres pairs

Avec Python, répondez à la question suivante en une seule commande. Combien y a-t-il de nombres pairs dans l'intervalle `[2, 10000]` inclus ?

Chapitre 5

Boucles et comparaisons

5.1 Boucles for

5.1.1 Principe

En programmation, on est souvent amené à répéter plusieurs fois une instruction. Incontournables à tout langage de programmation, les boucles vont nous aider à réaliser cette tâche de manière compacte. Imaginez par exemple que vous souhaitiez afficher les éléments d'une liste les uns après les autres. Dans l'état actuel de vos connaissances, il faudrait taper quelque chose du style :

```
1 animaux = ['girafe', 'tigre', 'singe', 'souris']
2 print(animaux[0])
3 print(animaux[1])
4 print(animaux[2])
5 print(animaux[3])
```

Si votre liste ne contient que quatre éléments, ceci est encore faisable mais imaginez qu'elle en contienne 100 voire 1000 ! Pour remédier à cela, il faut utiliser les boucles. Regardez l'exemple suivant :

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

Commentons en détails ce qu'il s'est passé dans cet exemple :

La variable `animal` est appelée **variable d'itération**, elle prend successivement les différentes valeurs de la liste `animaux` à chaque itération de la boucle. On verra en section 5.1.3 que l'on peut choisir le nom que l'on veut pour cette variable. Celle-ci est créée par Python la première fois que la ligne contenant le `for` est exécutée (si elle existait déjà son contenu serait écrasé). Une fois la boucle terminée, cette variable d'itération `animal` ne sera pas détruite et ainsi contiendra la dernière valeur de la liste `animaux` (ici la chaîne `souris`).

Notez bien les types de variables : `animaux` est une **liste** sur laquelle on itère, et `animal` est une **chaîne de caractères** car chaque élément de la liste est une chaîne de caractères. Nous verrons plus loin que la variable d'itération peut être de n'importe quel type selon la liste parcourue. En Python, une

boucle itère toujours sur un objet dit **séquentiel** (c'est à dire un objet constitué d'autres objets) tel qu'une liste. Nous verrons aussi plus tard d'autres objets séquentiels sur lesquels on peut itérer dans une boucle.

D'ores et déjà, remarquez avec attention le **signe deux-points** : à la fin de la ligne `for`. Cela signifie que la boucle `for` attend un **bloc d'instructions**, en l'occurrence toutes les instructions que Python répètera à chaque itération de la boucle. On appelle ce bloc d'instructions le **corps de la boucle**. Comment indique-t-on à Python où ce bloc commence et se termine? Cela est signalé uniquement par l'**indentation**, c'est-à-dire le décalage de la (ou des) ligne(s) du bloc d'instructions. Dans l'exemple suivant, le corps de la boucle contient deux instructions `print(animal*2)` et `print(animal)` car elles sont indentées par rapport à la ligne `for` :

```
1 for animal in animaux:
2     print(animal)
3     print(animal*2)
4 print("C'est fini")
```

La ligne `print("C'est fini")` ne fait pas partie du corps de la boucle car elle est au même niveau que le `for` (c'est-à-dire non indentée par rapport au `for`). Notez également que chaque instruction du corps de la boucle doit être indentée de la même manière (ici 4 espaces).

Outre une meilleure lisibilité, les `:` et l'**indentation** sont formellement requis en Python. Même si on peut indenter comme on veut en Python (plusieurs espaces ou plusieurs tabulations, mais pas une combinaison des deux), les développeurs recommandent l'utilisation de 4 espaces¹. Faites en sorte de configurer votre éditeur favori de façon à écrire 4 espaces lorsque vous tapez sur la touche Tab.

Si on oublie l'indentation, Python vous renvoie un message d'erreur :

```
1 >>> for animal in animaux:
2     ... print(animal)
3     File "<stdin>", line 2
4     print(animal)
5         ^
6 IndentationError: expected an indented block
```

Dans les exemples ci-dessus, nous avons exécuté une boucle en itérant directement sur une liste. Une tranche d'une liste étant elle même une liste, on peut également itérer dessus :

```
1 >>> animaux = ['girafe', 'tigre', 'singe', 'souris']
2 >>> for animal in animaux[1:3]:
3     ...     print(animal)
4     ...
5 tigre
6 singe
```

On a vu que les boucles `for` pouvaient utiliser une liste contenant des chaînes de caractères, mais elles peuvent tout aussi bien utiliser des listes contenant des entiers (ou n'importe quel type de variable finalement).

```
1 >>> for i in [1,2,3]:
2     ...     print(i)
3     ...
4 1
5 2
6 3
```

1. <https://www.python.org/dev/peps/pep-0008/#id17>

5.1.2 Fonction range()

Python possède la fonction `range()` que nous avons rencontrée précédemment dans le chapitre sur les *Listes* et qui est aussi bien commode pour faire une boucle sur une liste d'entiers de manière automatique :

```

1 >>> for i in range(4):
2     ...     print(i)
3     ...
4     0
5     1
6     2
7     3

```

Dans cet exemple nous pouvons faire plusieurs remarques importantes :

Contrairement à la génération de liste avec `list(range(4))`, la fonction `range()` peut être utilisée telle quelle dans une boucle. Il n'est pas nécessaire de taper `for i in list(range(4)):`, même si cela fonctionnerait également.

Comment cela est-il possible ? Et bien `range()` est une fonction qui a été spécialement conçue pour cela², c'est à dire que l'on peut itérer directement dessus. Pour Python, il s'agit d'un nouveau type, par exemple dans `x = range(3)` la variable `x` est de type *range* (tout comme on avait les types *int*, *float*, *str* ou *list*) à utiliser spécialement avec les boucles.

L'instruction `list(range(4))` se contente de transformer un objet de type *range* en un objet de type *list* (si vous vous souvenez bien, il s'agit donc d'une fonction de *casting*, c'est à dire de conversion d'un type en un autre). Il n'y a aucun intérêt à utiliser dans une boucle la construction `for i in list(range(4)):`. C'est même contre-productif. En effet, `range()` se contente de stocker l'entier actuel, le pas pour passer à l'entier suivant, et le dernier entier à parcourir, ce qui revient à stocker seulement 3 entiers et ce quel que soit la longueur de la séquence, même avec un `range(1000000)`. Si on utilisait `list(range(1000000))`, Python construirait d'abord une liste de 1 million d'éléments dans la mémoire puis itérerait dessus, d'où une énorme perte de temps !

5.1.3 Nommage de la variable d'itération

Dans l'exemple précédent, nous avons choisi le nom `i` pour la variable d'itération. Ceci est un standard en informatique et indique en général qu'il s'agit d'un entier (le nom `i` vient sans doute du mot indice ou *index* en anglais). Nous vous conseillons de suivre cette convention afin d'éviter les confusions, si vous itérez sur les indices vous pouvez appeler la variable d'itération `i` (par exemple dans `for i in range(4):`).

Si vous itérez sur une liste comportant des chaînes de caractères, mettez un nom explicite pour la variable d'itération. Par exemple :

```

1 for prenom in ['Joe', 'Bill', 'John']:

```

5.1.4 Itération sur les indices

Revenons à notre liste `animaux`. Nous allons maintenant parcourir cette liste, mais cette fois par une itération sur ses indices :

2. <https://docs.python.org/3/library/stdtypes.html#typeseq-range>

```
1 >>> animaux = ['girafe', 'tigre', 'singe', 'souris']
2 >>> for i in range(4):
3     ...     print(animaux[i])
4     ...
5 girafe
6 tigre
7 singe
8 souris
```

La variable `i` prendra les valeurs successives 0, 1, 2 et 3 et on accèdera à chaque élément de la liste `animaux` par son indice (*i.e.* `animaux[i]`). Notez à nouveau le nom `i` de la variable d'itération car on itère sur les **indices**.

Quand utiliser l'une ou l'autre des 2 méthodes? La plus efficace est celle qui réalise **les itérations directement sur les éléments** :

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

Toutefois, il se peut qu'au cours d'une boucle vous ayez besoin des indices, auquel cas vous n'avez pas le choix de itérer sur les indices. Par exemple :

```
1 >>> animaux = ['girafe','tigre','singe','souris']
2 >>> for i in range(len(animaux)):
3     ...     print("L'animal {} est un(e) {}".format(i, animaux[i]))
4     ...
5 L'animal 0 est un(e) girafe
6 L'animal 1 est un(e) tigre
7 L'animal 2 est un(e) singe
8 L'animal 3 est un(e) souris
```

Python possède également la fonction `enumerate()` qui vous permet d'itérer sur les indices et les éléments eux-mêmes.

```
1 >>> animaux = ['girafe','tigre','singe','souris']
2 >>> for i, animal in enumerate(animaux):
3     ...     print("L'animal {} est un(e) {}".format(i, animal))
4     ...
5 L'animal 0 est un(e) girafe
6 L'animal 1 est un(e) tigre
7 L'animal 2 est un(e) singe
8 L'animal 3 est un(e) souris
```

5.2 Comparaisons

Avant de passer à une autre sorte de boucles (les boucles `while`), nous abordons tout de suite les **comparaisons**. Celles-ci seront reprises dans le chapitre sur les *Tests*.

Python est capable d'effectuer toute une série de comparaisons entre le contenu de deux variables, telles que :

Syntaxe Python	Signification
<code>==</code>	égal à
<code>!=</code>	différent de
<code>></code>	supérieur à
<code>>=</code>	supérieur ou égal à
<code><</code>	inférieur à
<code><=</code>	inférieur ou égal à

Observez l'exemple suivant avec des nombres entiers.

```

1 >>> x = 5
2 >>> x == 5
3 True
4 >>> x > 10
5 False
6 >>> x < 10
7 True

```

Python renvoie la valeur `True` si la comparaison est vraie et `False` si elle est fausse. `True` et `False` sont des booléens.

Faites bien attention à ne pas confondre l'**opérateur d'affectation** `=` qui affecte une valeur à une variable et l'**opérateur de comparaison** `==` qui compare les valeurs de deux variables.

Vous pouvez également effectuer des comparaisons sur des chaînes de caractères.

```

1 >>> animal = "tigre"
2 >>> animal == "tig"
3 False
4 >>> animal != "tig"
5 True
6 >>> animal == 'tigre'
7 True

```

Dans le cas des chaînes de caractères, *a priori* seuls les tests `==` et `!=` ont un sens. En fait, on peut aussi utiliser les opérateurs `<`, `>`, `<=` et `>=`. Dans ce cas l'ordre alphabétique est pris en compte, par exemple :

```

1 >>> "a" < "b"
2 True

```

"a" est *inférieur* à "b" car il est situé avant dans l'ordre alphabétique. En fait, c'est l'ordre ASCII³ des caractères qui est pris en compte (*i.e.* chaque caractère est affecté à un code numérique), on peut donc comparer aussi des caractères spéciaux (comme `#` ou `~`) entre eux. On peut aussi comparer des chaînes à plus d'un caractère.

3. http://fr.wikipedia.org/wiki/American_Standard_Code_for_Information_Interchange

```
1 >>> "ali" < "alo"
2 True
3 >>> "abb" < "ada"
4 True
```

Dans ce cas, Python compare caractère par caractère de la gauche vers la droite (le premier avec le premier, le deuxième avec le deuxième, etc). Dès qu'un caractère est différent entre l'une et l'autre des chaînes, il considère que la chaîne la plus petite est celle qui présente le caractère ayant le plus petit code ASCII (les caractères suivants de la chaîne sont ignorés dans la comparaison), comme dans l'exemple "abb" < "ada" ci-dessus.

5.3 Boucles while

Une autre alternative à l'instruction `for` couramment utilisée en informatique est la boucle `while`. Le principe est simple. Une série d'instructions est exécutée tant qu'une condition est vraie. Par exemple :

```
1 >>> i = 1
2 >>> while i <= 4:
3     ...     print(i)
4     ...     i = i + 1
5     ...
6     1
7     2
8     3
9     4
```

Remarquez qu'il est encore une fois nécessaire d'indenter le bloc d'instructions correspondant au corps de la boucle (ici 2 instructions).

Une boucle `while` nécessite généralement **trois éléments** pour fonctionner correctement :

- l'initialisation de la variable de test avant la boucle ;
- le test de la variable associé à l'instruction `while` ;
- la mise à jour de la variable de test dans le corps de la boucle.

Faites bien attention aux tests et à l'incréméntation que vous utilisez car une erreur mène souvent à des boucles infinies, c'est-à-dire qui ne s'arrêtent jamais. Vous pouvez néanmoins toujours stopper l'exécution d'un script Python à l'aide de la combinaison de touches Ctrl-C. Par exemple :

```
1 i = 0
2 while i < 10:
3     print("Le python c'est cool !")
```

Ici nous avons omis de mettre à jour la variable `i`, ainsi la boucle ne s'arrêtera jamais (sauf en pressant Ctrl-C) puisque la condition `i < 10` sera toujours vraie.

5.4 Exercices

Conseil : pour ces exercices, écrivez des scripts dans des fichiers, puis exécutez-les dans un *shell*.

5.4.1 Boucles de base

Soit la liste `['vache', 'souris', 'levure', 'bacterie']`. Affichez l'ensemble des éléments de cette liste (un élément par ligne) de trois manières différentes (deux avec `for` et une avec `while`).

5.4.2 Boucle et jours de la semaine

Constituez une liste `semaine` contenant les 7 jours de la semaine.

Écrivez une série d'instructions affichant les jours de la semaine (en utilisant une boucle `for`), ainsi qu'une autre série d'instructions affichant les jours du week-end (en utilisant une boucle `while`).

5.4.3 Table des 1

Avec une boucle, affichez les nombres de 1 à 10 sur une seule ligne.

5.4.4 Nombres pairs et impairs

Soit `impairs` la liste de nombres `[1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 21]`. Écrivez un programme qui, à partir de la liste `impairs`, construit une liste `pairs` dans laquelle tous les éléments de `impairs` sont incrémentés de 1.

5.4.5 Calcul de moyenne

Voici les notes d'un étudiant `[14, 9, 6, 8, 12]`. Calculez la moyenne de ces notes. Utilisez l'écriture formatée pour afficher la valeur de la moyenne avec deux décimales.

5.4.6 Produit de nombres consécutifs

Soit la liste `X` contenant les nombres entiers de 0 à 10. Calculez le produit des nombres consécutifs deux à deux de `X` en utilisant une boucle. Exemple pour les premières itérations :

```
1 0
2 2
3 6
4 12
```

5.4.7 Triangle

Écrivez un script qui dessine un triangle comme celui-ci :

```
1 *
2 **
3 ***
4 ****
5 *****
6 *****
7 *****
8 *****
```

```

9 *****
10 *****

```

5.4.8 Triangle inversé

Écrivez un script qui dessine un triangle comme celui-ci :

```

1 *****
2 *****
3 *****
4 *****
5 *****
6 *****
7 *****
8 *****
9 *****
10 *****

```

5.4.9 Triangle gauche

Écrivez un script qui dessine un triangle comme celui-ci :

```

1      *
2     **
3    ***
4   ****
5  *****
6 ******
7 *******
8 *******
9 *******
10 *******

```

5.4.10 Triangle isocèle

Écrivez un script qui dessine un triangle comme celui-ci :

```

1      *
2     ***
3    *****
4   *******
5  *********
6  *********
7  *********
8  *********
9  *********
10  *********

```

5.4.11 Parcours de matrice

Imaginons que l'on souhaite parcourir tous les éléments d'une matrice carrée, c'est-à-dire d'une matrice qui est constituée d'autant de colonnes que de lignes.

Écrivez un script qui parcourt chaque élément de la matrice et qui affiche le numéro de ligne et de colonne.

Pour une matrice 2x2, le schéma 5.1 vous indique comment parcourir une telle matrice. L'affichage attendu est :

```

1 ligne colonne
2   1     1
3   1     2
4   2     1
5   2     2

```

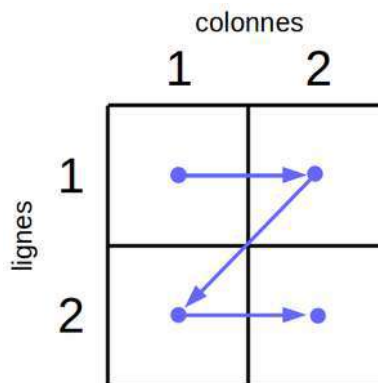
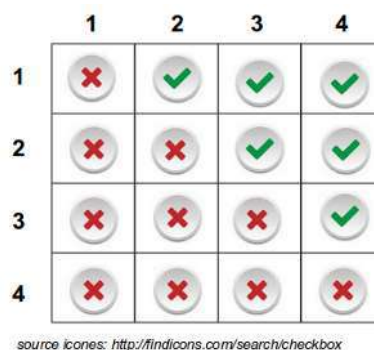


FIGURE 5.1 – Parcours d'une matrice

Attention à bien respecter l'alignement des chiffres qui doivent être justifiés à droite. Testez pour une matrice 3x3, 5x5, et 10x10.

5.4.12 Parcours de demi-matrice sans la diagonale (exercice ++)

En se basant sur le script précédent, on souhaite réaliser le parcours d'une demi-matrice carrée sans la diagonale. On peut noter que cela donne tous les couples possibles une seule fois (1 et 2 est équivalent à 2 et 1), en excluant par ailleurs chaque élément avec lui-même (1 et 1, 2 et 2, etc). Pour mieux comprendre ce qui est demandé, la figure 5.2 indique les cases à parcourir en vert :



source icones: <http://findicons.com/search/checkbox>

FIGURE 5.2 – Demi-matrice sans la diagonale (en vert)

Écrivez un script qui affiche le numéro de ligne et de colonne, puis la taille de la matrice $N \times N$ et le nombre total de cases parcourues. Par exemple pour une matrice 4×4 ($N=4$) :

```

1 ligne colonne
2   1   2
3   1   3
4   1   4
5   2   3
6   2   4
7   3   4
8 Pour une matrice 4x4, on a parcouru 6 cases

```

Testez votre script avec $N=3$, $N=4$ et $N=5$.

Concevez une seconde version à partir du script précédent, où cette fois on n’affiche plus tous les couples possibles mais simplement la valeur de N , et le nombre de cases parcourues. Affichez cela pour des valeurs de N allant de 1 à 20. Pouvez-vous trouver une formule générale reliant le nombre de cases parcourues à N ?

5.4.13 Sauts de puce

On imagine une puce qui se déplace aléatoirement sur une ligne, en avant ou en arrière, par pas de 1 ou -1. Par exemple, si elle est à l’emplacement 0, elle peut sauter à l’emplacement 1 ou -1 ; si elle est à l’emplacement 2, elle peut sauter à l’emplacement 3 ou 1, etc. Avec une boucle `while`, simuler le mouvement de cette puce de l’emplacement 0 à l’emplacement 5 (cf schéma suivant). Combien de sauts sont nécessaires à ce parcours ? Relancez plusieurs fois le programme. Trouvez-vous le même nombre de sauts à chaque exécution ? Comparez avec votre voisin.

```

1   départ   arrivée
2     |       |
3     v       v
4 ...o-X-o-o-o-o-X-o...
5   -1 0 1 2 3 4 5 6

```

Conseil : vous utiliserez l’instruction

```
random.choice([-1,1])
```

qui renvoie au hasard les valeurs -1 ou 1 avec la même probabilité. Avant d’utiliser cette commande vous devrez mettre en haut du script la ligne

```
import random
```

Nous reverrons la signification de cette syntaxe particulière dans le chapitre sur les modules.

5.4.14 Suite de Fibonacci (exercice +++)

La suite de Fibonacci est une suite mathématique qui porte le nom de Leonardo Fibonacci, un mathématicien italien du XIII^e siècle. Initialement, cette suite a été conçue pour décrire la croissance d’une population de lapins mais elle peut également être utilisée pour décrire certains motifs géométriques retrouvés dans la nature (coquillages, fleurs de tournesol...).

Par définition, les deux premiers termes de la suite de Fibonacci sont 0 et 1. Ensuite, le terme au rang n est la somme des nombres aux rangs $n - 1$ et $n - 2$. Par exemple, les 10 premiers termes de la suite de Fibonacci sont 0, 1, 1, 2, 3, 5, 8, 13, 21, 34.

Écrivez un script qui construit la liste des 20 premiers termes de la suite de Fibonacci puis l'affiche.

Améliorez ce script en imprimant à chaque itération le rapport entre l'élément n et l'élément $n - 1$. Ce rapport tend-il vers une constante ? Si oui, laquelle ?

Chapitre 6

Tests

6.1 Définition

Les **tests** sont un élément essentiel à tout langage informatique si on veut lui donner un peu de complexité car ils permettent à l'ordinateur de prendre des décisions si telle ou telle condition est vraie ou fausse. Pour cela, Python utilise l'instruction `if` ainsi qu'une comparaison que nous avons abordée au chapitre précédent. Voici un premier exemple :

```
1 >>> x = 2
2 >>> if x == 2:
3 ...     print("Le test est vrai !")
4 ...
5 Le test est vrai !
```

et un second :

```
1 >>> x = "souris"
2 >>> if x == "tigre":
3 ...     print("Le test est vrai !")
4 ...
```

Il y a plusieurs remarques à faire concernant ces deux exemples :

- Dans le premier exemple, le test étant vrai, l'instruction `print("Le test est vrai !")` est exécutée. Dans le second exemple, le test est faux et rien n'est affiché.
- Les blocs d'instruction dans les tests doivent forcément être indentés comme les boucles `for` et `while`. L'indentation indique la portée des instructions à exécuter si le test est vrai.
- L'instruction `if` se termine comme les instructions `for` et `while` par le caractère `:`.

6.2 Tests à plusieurs cas

Parfois, il est pratique de tester si la condition est vraie ou si elle est fausse dans une même instruction `if`. Plutôt que d'utiliser deux instructions `if`, on peut se servir des instructions `if` et `else` :

```
1 >>> x = 2
2 >>> if x == 2:
3 ...     print("Le test est vrai !")
4 ... else:
5 ...     print("Le test est faux !")
```

```

6 ...
7 Le test est vrai !
8 >>> x = 3
9 >>> if x == 2:
10 ...     print("Le test est vrai !")
11 ... else:
12 ...     print("Le test est faux !")
13 ...
14 Le test est faux !

```

On peut utiliser une série de tests dans la même instruction `if`, notamment pour tester plusieurs valeurs d'une même variable. Par exemple, on se propose de tirer au sort une base d'ADN puis d'afficher le nom de cette dernière. Dans le code suivant, nous utilisons l'instruction `random.choice(liste)` qui renvoie un élément choisi au hasard dans une `liste`. L'instruction `import random` sera vue plus tard, admettez pour le moment qu'elle est nécessaire.

```

1 >>> import random
2 >>> base = random.choice(["a", "t", "c", "g"])
3 >>> if base == "a":
4 ...     print("choix d'une adénine")
5 ... elif base == "t":
6 ...     print("choix d'une thymine")
7 ... elif base == "c":
8 ...     print("choix d'une cytosine")
9 ... elif base == "g":
10 ...    print("choix d'une guanine")
11 ...
12 choix d'une cytosine

```

Dans cet exemple, Python teste la première condition, puis, si et seulement si elle est fautive, teste la deuxième et ainsi de suite... Le code correspondant à la première condition vérifiée est exécuté puis Python sort du `if`.

Remarque : De nouveau, faites bien attention à l'indentation dans ces deux derniers exemples! Vous devez être très rigoureux sur ce point. Pour vous en convaincre, exécutez ces deux scripts dans l'interpréteur Python :

Script 1

```

1 nombres = [4, 5, 6]
2 for nb in nombres:
3     if nb == 5:
4         print("Le test est vrai")
5         print("car la variable nb vaut {}".format(nb))

```

Script 2

```

1 nombres = [4, 5, 6]
2 for nb in nombres:
3     if nb == 5:
4         print("Le test est vrai")
5     print("car la variable nb vaut {}".format(nb))

```

Comment expliquez-vous ce résultat? Observez bien l'indentation de la dernière ligne.

6.3 Tests multiples

Les tests multiples permettent de tester plusieurs conditions en même temps en utilisant des opérateurs booléens. Les deux opérateurs les plus couramment utilisés sont le **OU** et le **ET**. Voici un petit rappel du mode de fonctionnement de l'opérateurs **OU** :

Condition 1	Opérateur	Condition 2	Résultat
Vrai	OU	Vrai	Vrai
Vrai	OU	Faux	Vrai
Faux	OU	Vrai	Vrai
Faux	OU	Faux	Faux

et de l'opérateur **ET** :

Condition 1	Opérateur	Condition 2	Résultat
Vrai	ET	Vrai	Vrai
Vrai	ET	Faux	Faux
Faux	ET	Vrai	Faux
Faux	ET	Faux	Faux

En Python, on utilise le mot réservé `and` pour l'opérateur **ET** et le mot réservé `or` pour l'opérateur **OU**. Respectez bien la casse des opérateurs `and` et `or` qui, en Python, s'écrivent en minuscule. En voici un exemple d'utilisation :

```

1 >>> x = 2
2 >>> y = 2
3 >>> if x == 2 and y == 2:
4 ...     print("le test est vrai")
5 ...
6 le test est vrai

```

Notez que le même résultat serait obtenu en utilisant deux instructions `if` imbriquées :

```

1 >>> x = 2
2 >>> y = 2
3 >>> if x == 2:
4 ...     if y == 2:
5 ...         print("le test est vrai")
6 ...
7 le test est vrai

```

Vous pouvez aussi tester directement l'effet de ces opérateurs à l'aide de `True` et `False` (attention à respecter la casse).

```

1 >>> True or False
2 True

```

Enfin, on peut utiliser l'opérateur logique de négation `not` qui inverse le résultat d'une condition :

```

1 >>> not True
2 False
3 >>> not False

```

```

4 True
5 >>> not (True and True)
6 False

```

6.4 Instructions break et continue

Ces deux instructions permet de modifier le comportement d'une boucle (**for** ou **while**) avec un test.

L'instruction **break** stoppe la boucle.

```

1 >>> for i in range(5):
2 ...     if i > 2:
3 ...         break
4 ...     print(i)
5 ...
6 0
7 1
8 2

```

L'instruction **continue** saute à l'itération suivante.

```

1 >>> for i in range(5):
2 ...     if i == 2:
3 ...         continue
4 ...     print(i)
5 ...
6 0
7 1
8 3
9 4

```

6.5 Tests de valeur sur des réels

Lorsque l'on souhaite tester la valeur d'une variable *float*, le premier réflexe serait d'utiliser l'opérateur d'égalité comme :

```

1 >>> 1/10 == 0.1
2 True

```

Toutefois nous vous le conseillons formellement. Pourquoi? Python stocke les *float* sous forme de nombres flottants (d'où leur nom!), et cela mène à certaines limitations¹. Regardez l'exemple suivant :

```

1 >>> (3 - 2.7) == 0.3
2 False
3 >>> 3 - 2.7
4 0.29999999999999998

```

Nous voyons que le résultat de l'opération $3 - 2.7$ n'est pas exactement 0.3 d'où le **False**. Pour éviter ces problèmes nous conseillons de toujours encadrer un *float* avec un niveau de précision *delta*, par exemple :

1. <https://docs.python.org/3.6/tutorial/floatpoint.html>

```

1 >>> delta = 0.0001
2 >>> 0.3 - delta < 3.0 - 2.7 < 0.3 + delta
3 True

```

Ici on teste si $3 - 2.7$ est compris entre $0.3 \pm \text{delta}$.

6.6 Exercices

Conseil : pour ces exercices, écrivez des scripts dans des fichiers, puis exécutez-les dans un *shell*.

6.6.1 Jours de la semaine

Constituez une liste `semaine` contenant les sept jours de la semaine.

En utilisant une boucle, écrivez chaque jour de la semaine ainsi que les messages suivants :

- Au travail s'il s'agit du lundi au jeudi;
- Chouette c'est vendredi s'il s'agit du vendredi;
- Repos ce week-end s'il s'agit du week-end.

Ces messages ne sont que des suggestions, vous pouvez laisser libre cours votre imagination.

6.6.2 Séquence complémentaire d'un brin d'ADN

La liste ci-dessous représente la séquence d'un brin d'ADN :

```
["A", "C", "G", "T", "T", "A", "G", "C", "T", "A", "A", "C", "G"]
```

Écrivez un script qui transforme cette séquence en sa séquence complémentaire.

Rappel : la séquence complémentaire s'obtient en remplaçant A par T, T par A, C par G et G par C.

6.6.3 Minimum d'une liste

La fonction `min()` de Python, renvoie l'élément le plus petit d'une liste constituée de valeurs numériques ou de chaînes de caractères. Sans utiliser cette fonction, écrivez un script qui détermine le plus petit élément de la liste `[8, 4, 6, 1, 5]`.

6.6.4 Fréquence des acides aminés

La liste ci-dessous représente une séquence d'acides aminés :

```
["A", "R", "A", "W", "W", "A", "W", "A", "R", "W", "W", "R", "A", "G", "A", "R"]
```

Calculez la fréquence des acides aminés alanine (A), arginine (R), tryptophane (W) et glycine (G) dans cette séquence.

6.6.5 Notes et mention d'un étudiant

Voici les notes d'un étudiant : 14, 9, 13, 15 et 12. Écrivez un script qui affiche la note maximum (fonction `max()`), la note minimum (fonction `min()`) et qui calcule la moyenne.

Affichez la valeur de la moyenne avec deux décimales. Affichez aussi la mention obtenue sachant que la mention est passable si la moyenne est entre 10 inclus et 12 exclus, assez-bien entre 12 inclus et 14 exclus et bien au-delà de 14.

6.6.6 Nombres pairs

Construisez une boucle qui parcourt les nombres de 0 à 20 et qui affiche les nombres pairs inférieurs ou égaux à 10 d'une part, et les nombres impairs strictement supérieur à 10 d'autre part.

Pour cet exercice, vous pourrez utiliser l'opérateur modulo `%` qui renvoie le reste de la division entière entre deux nombres et dont voici quelques exemples d'utilisation :

```

1 >>> 4 % 3
2 1
3 >>> 5 % 3
4 2
5 >>> 4 % 2
6 0
7 >>> 5 % 2
8 1
9 >>> 6 % 2
10 0
11 >>> 7 % 2
12 1

```

Ainsi, vous remarquerez qu'un nombre est pair lorsque le reste de sa division entière par 2 est nul.

6.6.7 L'énigme du père Fouras

A Fort Boyard, le père Fouras nous pose l'énigme suivante :

Pour ouvrir le coffre où se trouve la clé, trouve la combinaison à trois chiffres sachant que :

- *Le nombre est inférieur à 200.*
- *Deux de ses chiffres sont identiques.*
- *La somme de ses chiffres est égale à 5.*
- *C'est un nombre pair.*

On se propose d'utiliser une méthode dite *brute force*, c'est à dire d'utiliser une boucle et à chaque itération on teste les 4 conditions. Avez-vous réussi à trouver la même solution en raisonnant ?

6.6.8 Conjecture de Syracuse (exercice +++)

La conjecture de Syracuse² est une conjecture mathématique qui reste improuvée à ce jour et qui est définie de la manière suivante.

2. http://fr.wikipedia.org/wiki/Conjecture_de_Syracuse

Soit un entier positif n . Si n est pair, alors le diviser par 2. Si il est impair, alors le multiplier par 3 et lui ajouter 1. En répétant cette procédure, la suite de nombres atteint la valeur 1 puis se prolonge indéfiniment par une suite de trois valeurs triviales appelée cycle trivial.

Jusqu'à présent, la conjecture de Syracuse, selon laquelle depuis n'importe quel entier positif la suite de Syracuse atteint 1, n'a pas été mise en défaut.

Par exemple, les premiers éléments de la suite de Syracuse si on prend comme point de départ 10 sont : 10, 5, 16, 8, 4, 2, 1...

Écrivez un script qui, partant d'un entier positif n (par exemple 10 ou 20), crée une liste des nombres de la suite de Syracuse. Avec différents points de départ (c'est-à-dire avec différentes valeurs de n), la conjecture de Syracuse est-elle toujours vérifiée ? Quels sont les nombres qui constituent le cycle trivial ?

Remarques

1. Pour cet exercice, vous avez besoin de faire un nombre d'itérations inconnu pour que la suite de Syracuse atteigne le chiffre 1 puis entame son cycle trivial. Vous pourrez tester votre algorithme avec un nombre arbitraire d'itérations, typiquement 20 ou 100, suivant votre nombre n de départ.
2. Un nombre est pair lorsque le reste de sa division entière (opérateur modulo %) par 2 est nul.

6.6.9 Attribution simple de la structure secondaire des résidus d'une protéine (exercice +++)

Les angles dièdres phi/psi d'une hélice alpha parfaite ont une valeur de -57 degrés et -47 degrés respectivement. Bien sûr, il est très rare que l'on trouve ces valeurs parfaites dans une protéine, par conséquent il est couramment accepté de tolérer une déviation de +/- 30 degrés sur celles-ci.

Vous trouverez ci-dessous une liste de listes contenant les valeurs des angles phi/psi des résidues de la première hélice de la protéine 1TFE³. En utilisant cette liste, écrivez un programme qui teste, pour chaque résidu, s'il est ou pas en hélice.

```

1 [[48.6,53.4], [-124.9,156.7], [-66.2,-30.8], [-58.8,-43.1], \
2 [-73.9,-40.6], [-53.7,-37.5], [-80.6,-16.0], [-68.5,135.0], \
3 [-64.9,-23.5], [-66.9,-45.5], [-69.6,-41.0], [-62.7,-37.5], \
4 [-68.2,-38.3], [-61.2,-49.1], [-59.7,-41.1], [-63.2,-48.5], \
5 [-65.5,-38.5], [-64.1,-40.7], [-63.6,-40.8], [-66.4,-44.5], \
6 [-56.0,-52.5], [-55.4,-44.6], [-58.6,-44.0], [-77.5,-39.1], \
7 [-91.7,-11.9], [48.6,53.4]]

```

6.6.10 Détermination des nombres premiers inférieurs à 100 (exercice +++)

Voici un extrait de l'article sur les nombres premiers tiré de l'encyclopédie en ligne wikipédia⁴.

Un nombre premier est un entier naturel qui admet exactement deux diviseurs distincts entiers et positifs (qui sont alors 1 et lui-même). Cette définition exclut 1, qui n'a qu'un seul diviseur entier positif. Par opposition, un nombre non nul produit de deux nombres entiers différents de 1 est dit composé. Par exemple $6 = 2 \times 3$ est composé, tout comme $21 = 3 \times 7$, mais 11 est premier car 1 et 11 sont les seuls diviseurs de 11. Les nombres 0 et 1 ne sont ni premiers ni composés.

3. <http://www.rcsb.org/pdb/explore.do?structureId=1TFE>

4. http://fr.wikipedia.org/wiki/Nombre_premier

Déterminez les nombres premiers inférieurs à 100. Combien y a-t-il de nombres premiers entre 0 et 100 ? Pour vous aider, nous vous proposons deux méthodes.

Méthode 1 (peu optimale mais assez intuitive)

Pour chaque nombre de 2 à 100, calculez le reste de la division entière (avec l'opérateur modulo %) depuis 1 jusqu'à lui-même. Si c'est un nombre premier, il aura exactement deux nombres pour lesquels le reste de la division entière est égal à 0 (1 et lui-même). Si ce n'est pas un nombre premier, il aura plus de deux nombres pour lesquels le reste de la division entière est égal à 0.

Méthode 2 (plus optimale et plus rapide, mais un peu plus compliquée)

Vous pouvez parcourir tous les nombres de 2 à 100 et vérifier si ceux-ci sont composés, c'est-à-dire qu'ils sont le produit de deux nombres premiers. Pratiquement, cela consiste à vérifier que le reste de la division entière (opérateur modulo %) entre le nombre considéré et n'importe quel nombre premier est nul. Le cas échéant, ce nombre n'est pas premier.

6.6.11 Recherche d'un nombre par dichotomie (exercice +++)

La recherche par dichotomie⁵ est une méthode qui consiste à diviser (en générale en parties égales) un problème pour en trouver la solution. À titre d'exemple, voici une discussion entre Pierre et Patrick dans laquelle Pierre essaie de deviner le nombre (compris entre 1 et 100) auquel Patrick a pensé.

- [Patrick] “C'est bon, j'ai pensé à un nombre entre 1 et 100.”
- [Pierre] “OK, je vais essayer de le deviner. Est-ce que ton nombre est plus petit ou plus grand que 50 ?”
- [Patrick] “Plus grand.”
- [Pierre] “Est-ce que ton nombre est plus petit, plus grand ou égal à 75 ?”
- [Patrick] “Plus grand.”
- [Pierre] “Est-ce que ton nombre est plus petit, plus grand ou égal à 87 ?”
- [Patrick] “Plus petit.”
- [Pierre] “Est-ce que ton nombre est plus petit, plus grand ou égal à 81 ?”
- [Patrick] “Plus petit.”
- [Pierre] “Est-ce que ton nombre est plus petit, plus grand ou égal à 78 ?”
- [Patrick] “Plus grand.”
- [Pierre] “Est-ce que ton nombre est plus petit, plus grand ou égal à 79 ?”
- [Patrick] “Egal. C'est le nombre auquel j'avais pensé. Bravo !”

Pour arriver rapidement à deviner le nombre, l'astuce consiste à prendre à chaque fois la moitié de l'intervalle dans lequel se trouve le nombre. Voici le détail des différentes étapes :

1. le nombre se trouve entre 1 et 100, on propose 50 ($100 / 2$).
2. le nombre se trouve entre 50 et 100, on propose 75 ($50 + (100-50)/2$).
3. le nombre se trouve entre 75 et 100, on propose 87 ($75 + (100-75)/2$).
4. le nombre se trouve entre 75 et 87, on propose 81 ($75 + (87-75)/2$).
5. le nombre se trouve entre 75 et 81, on propose 78 ($75 + (81-75)/2$).
6. le nombre se trouve entre 78 et 81, on propose 79 ($78 + (81-78)/2$).

Écrivez un script qui reproduit ce jeu de devinettes. Vous pensez à un nombre entre 1 et 100 et l'ordinateur essaie de le deviner par dichotomie en vous posant des questions.

Vous utiliserez la fonction `input()` pour interagir avec l'utilisateur. Voici un exemple de son fonctionnement :

5. <https://fr.wikipedia.org/wiki/Dichotomie>

```
1 >>> lettre = input("Entrez une lettre : ")
2 Entrez une lettre : P
3 >>> print(lettre)
4 P
```

Pour vous guider, voici ce que donnerait le programme avec la conversation précédente :

```
1 Pensez à un nombre entre 1 et 100.
2 Est-ce votre nombre est plus grand, plus petit ou égal à 50 ? [+/-/=] +
3 Est-ce votre nombre est plus grand, plus petit ou égal à 75 ? [+/-/=] +
4 Est-ce votre nombre est plus grand, plus petit ou égal à 87 ? [+/-/=] -
5 Est-ce votre nombre est plus grand, plus petit ou égal à 81 ? [+/-/=] -
6 Est-ce votre nombre est plus grand, plus petit ou égal à 78 ? [+/-/=] +
7 Est-ce votre nombre est plus grand, plus petit ou égal à 79 ? [+/-/=] =
8 J'ai trouvé en 6 questions !
```

Les caractères [+/-/=] indiquent à l'utilisateur comment il doit interagir avec l'ordinateur, c'est-à-dire entrer soit le caractère + si le nombre choisi est plus grand que le nombre proposé par l'ordinateur, soit le caractère - si le nombre choisi est plus petit que le nombre proposé par l'ordinateur, soit le caractère = si le nombre choisi est celui proposé par l'ordinateur (en appuyant ensuite sur la touche *Entrée*).

Chapitre 7

Fichiers

7.1 Lecture dans un fichier

Dans la plupart des travaux de programmation, on doit lire ou écrire dans un fichier. Python possède pour cela tout un tas d'outils qui vous simplifient la vie. Avant de passer à un exemple concret, créez un fichier dans un éditeur de texte que vous enregistrerez dans votre répertoire avec le nom `zoo.txt`, par exemple :

```
1 girafe
2 tigre
3 singe
4 souris
```

Ensuite, testez le code suivant :

```
1 >>> filin = open('zoo.txt', 'r')
2 >>> filin
3 <_io.TextIOWrapper name='zoo.txt' mode='r' encoding='UTF-8'>
4 >>> filin.readlines()
5 ['girafe\n', 'tigre\n', 'singe\n', 'souris\n']
6 >>> filin.close()
7 >>> f.readlines()
8 Traceback (most recent call last):
9   File "<stdin>", line 1, in <module>
10 ValueError: I/O operation on closed file.
```

Il y a plusieurs commentaires à faire sur cet exemple :

- La première commande ouvre le fichier `zoo.txt` en lecture seule (ceci est indiqué avec la lettre `r`). Remarquez que le fichier n'est pas encore lu, mais simplement ouvert (*un peu comme lorsqu'on ouvre un livre, mais qu'on ne l'a pas encore lu*). Le curseur de lecture est prêt à lire le premier caractère du fichier. Lorsqu'on affiche la valeur de la variable `filin`, vous voyez que Python la considère comme un objet de type fichier ouvert : `<_io.TextIOWrapper name='zoo.txt' mode='r' encoding='UTF-8'>`.
- Ensuite, nous rencontrons à nouveau la syntaxe `objet.méthode()` (cf. Chapitre 3) : ici la méthode `readlines()` agit sur l'objet `filin` en déplaçant le curseur de lecture du début à la fin du fichier, puis elle renvoie une liste contenant les lignes du fichier (*dans notre analogie avec un livre, ceci correspondrait à lire toutes les lignes du livre*).
- Enfin, on applique la méthode `close()` sur l'objet `filin`, ce qui, vous vous en doutez, ferme le fichier (*ceci correspondrait à fermer le livre*). On pourra remarquer que la méthode `.close()`

ne renvoie rien mais modifie l'état de l'objet `filin` en fichier fermé. Ainsi, si on essaie de lire à nouveau les lignes du fichier, Python renvoie une erreur car il ne peut pas lire un fichier fermé.

Voici maintenant un exemple complet de lecture d'un fichier avec Python.

```
1 >>> filin = open('zoo.txt', 'r')
2 >>> lignes = filin.readlines()
3 >>> lignes
4 ['girafe\n', 'tigre\n', 'singe\n', 'souris\n']
5 >>> for ligne in lignes:
6 ...     print(ligne)
7 ...
8 girafe
9
10 tigre
11
12 singe
13
14 souris
15
16 >>> filin.close()
```

Vous voyez qu'en cinq lignes de code, vous avez lu et parcouru le fichier.

Remarques :

- Notez que la liste `lignes` contient le caractère `\n` à la fin de chacun de ses éléments. Ceci correspond au saut à la ligne de chacune d'entre elles (ceci est codé par un caractère spécial que l'on représente par `\n`). Vous pourrez parfois rencontrer également la notation octale `\012`.
- Remarquez aussi que lorsque l'on affiche les différentes lignes du fichier à l'aide de la boucle `for` et de l'instruction `print()`, Python saute à chaque fois une ligne.

7.1.1 Méthode `read()`

Il existe d'autres méthodes que `readlines()` pour lire (et manipuler) un fichier. Par exemple, la méthode `read()` lit tout le contenu d'un fichier et renvoie une chaîne de caractères unique.

```
1 >>> filin = open('zoo.txt', 'r')
2 >>> filin.read()
3 'girafe\ntigre\nsinge\nsouris\n'
4 >>> filin.close()
```

7.1.2 Méthode `readline()`

La méthode `readline()` (sans `s` à la fin) lit une ligne d'un fichier et la renvoie sous forme d'une chaîne de caractères. À chaque nouvel appel de `readline()`, la ligne suivante est renvoyée. Associée à la boucle `while`, cette méthode permet de lire un fichier ligne par ligne.

```
1 >>> filin = open('zoo.txt', 'r')
2 >>> ligne = filin.readline()
3 >>> while ligne != "":
4 ...     print(ligne)
5 ...     ligne = filin.readline()
6 ...
```

```
7 girafe
8
9 tigre
10
11 singe
12
13 souris
14
15 >>> filin.close()
```

7.1.3 Méthodes seek() et tell()

Les méthodes `seek()` et `tell()` permettent respectivement de se déplacer au nième caractère (plus exactement au nième octet) d'un fichier et d'afficher où en est la lecture du fichier, c'est-à-dire quel caractère (ou octet) est en train d'être lu.

```
1 >>> filin = open('zoo.txt', 'r')
2 >>> filin.readline()
3 'girafe\n'
4 >>> filin.tell()
5 7
6 >>> filin.seek(0)
7 0
8 >>> filin.tell()
9 0
10 >>> filin.readline()
11 'girafe\n'
12 >>> filin.close()
```

On remarque qu'à l'ouverture d'un fichier, le tout premier caractère est considéré comme le caractère 0 (tout comme le premier élément d'une liste). La méthode `seek()` permet facilement de remonter au début du fichier lorsque l'on est arrivé à la fin ou lorsqu'on en a lu une partie. En passant, `seek()` confirme le numéro d'octet auquel il s'est déplacé en l'affichant à l'écran.

7.1.4 Itérations directement sur le fichier

Python essaie de vous faciliter la vie au maximum. Voici un moyen à la fois simple et élégant de parcourir un fichier.

```
1 >>> filin = open('zoo.txt', 'r')
2 >>> for ligne in filin:
3     ...     print(ligne)
4     ...
5 girafe
6
7 tigre
8
9 singe
10
11 souris
12
13 >>> filin.close()
```

La boucle `for` va demander à Python d'aller lire le fichier ligne par ligne. Privilégiez cette méthode par la suite.

7.2 Écriture dans un fichier

Écrire dans un fichier est aussi simple que de le lire. Voyez l'exemple suivant :

```
1 >>> animaux2 = ['poisson', 'abeille', 'chat']
2 >>> filout = open('zoo2.txt', 'w')
3 >>> for animal in animaux2:
4 ...     filout.write(animal)
5 ...
6 7
7 7
8 4
9 >>> filout.close()
```

Le contenu du fichier `zoo2.txt` est `poissonabeillechat`.

Quelques commentaires sur cet exemple :

- Après avoir initialisé la liste `animaux2`, nous avons ouvert un fichier mais cette fois-ci en mode écriture (avec le caractère `w`).
- Ensuite, on a balayé cette liste à l'aide d'une boucle. À chaque itération, nous avons écrit chaque élément de la liste dans le fichier. Remarquez à nouveau la méthode `write()` qui s'applique sur l'objet `filout`. Notez qu'à chaque utilisation de la méthode `write()`, celle-ci nous affiche le nombre d'octets écrits dans le fichier. Ceci est valable uniquement dans l'interpréteur, si vous créez un programme avec les mêmes lignes de code, ces valeurs ne s'afficheront pas à l'écran.
- Enfin, on a fermé le fichier avec la méthode `close()`.

Vous voyez qu'il est extrêmement simple en Python de lire ou d'écrire dans un fichier.

Remarque : si votre programme produit uniquement du texte, vous pouvez l'écrire sur la sortie standard (avec l'instruction `print()`). L'avantage est que dans ce cas l'utilisateur peut bénéficier de toute les potentialités d'Unix (redirection, tri, *parsing*...). S'il veut écrire le résultat du programme dans un fichier, il pourra toujours le faire en redirigeant la sortie.

7.3 Méthode optimisée d'ouverture et de fermeture de fichier

Depuis la version 2.5, Python introduit le mot-clé `with` qui permet d'ouvrir et de fermer un fichier de manière commode. Si pour une raison ou une autre l'ouverture conduit à une erreur (problème de droits, etc), l'utilisation de `with` garantit la bonne fermeture du fichier (ce qui n'est pas le cas avec l'utilisation de la méthode `open()` invoquée telle quelle). Voici un exemple :

```
1 >>> with open('zoo.txt', 'r') as filin:
2 ...     for ligne in filin:
3 ...         print(ligne)
4 ...
5 girafe
6
7 tigre
8
9 singe
```

```
10
11 souris
12
13 >>>
```

Vous remarquez que `with` introduit un bloc d'indentation. C'est à l'intérieur de ce bloc que nous effectuons toutes les opérations sur le fichier. Une fois sorti, Python fermera **automatiquement** le fichier. Vous n'avez donc plus besoin d'invoquer la méthode `close()`.

La commande `with` permet même d'ouvrir deux fichiers en même temps ! Voyez l'exemple suivant :

```
1 with open('zoo.txt', 'r') as fichier1, open('zoo2.txt', 'w') as fichier2:
2     for ligne in fichier1:
3         fichier2.write('* ' + ligne)
```

Si le fichier `zoo.txt` contient le texte suivant :

```
1 souris
2 girafe
3 lion
4 singe
```

alors le contenu de `zoo2.txt` sera :

```
1 * souris
2 * girafe
3 * lion
4 * singe
```

Dans cet exemple, `with` permet une notation très compacte en s'affranchissant de deux `.close()` en même temps.

Pour ceux qui veulent approfondir, la commande `with`¹ est plus générale et utilisable dans d'autres contextes (méthode compacte pour gérer les exceptions).

7.4 Note sur les retours chariots sous Unix et sous Windows

Conseil : si vous êtes débutant, vous pouvez sauter cette section.

On a vu plus haut que le caractère spécial `\n` correspondait à un retour à la ligne. C'est le standard en Unix.

Toutefois, Windows utilise deux caractères spéciaux pour le retour à la ligne : `\r` correspondant à un retour chariot (hérité des machines à écrire !) et `\n` comme en Unix.

Si vous avez commencé en Python 2, vous aurez peut-être remarqué que selon les versions la lecture de fichier supprimait parfois les `\r` et d'autres fois les laissait. Heureusement, la fonction `open()` dans Python 3² gère tout ça automatiquement et renvoie uniquement des sauts de ligne sous forme de `\n` (même si le fichier a été conçu sous Windows et qu'il contient initialement des `\r`).

1. https://docs.python.org/3/reference/compound_stmts.html#the-with-statement

2. <https://docs.python.org/3.6/library/functions.html#open>

7.5 Importance des conversions de types avec les fichiers

Vous avez noté que les méthodes qui lisent un fichier (par exemple `readlines()`) vous renvoient systématiquement des chaînes de caractères. De même, pour écrire dans un fichier il faut passer une chaîne de caractères à la méthode `.write()`. Pour ce genre de problème, il faudra utiliser les fonctions de conversions de types vues au chapitre 2 `int()`, `float()` et `str()`. Ces conversions sont essentielles lorsqu'on lit ou écrit des nombres dans un fichier. En effet, les nombres dans un fichier sont considérés comme du texte par la fonction `readlines()`, par conséquent il faut les convertir si on veut effectuer des opérations numériques dessus.

7.6 Exercices

Conseil : pour tous ces exercices, écrivez des scripts dans des fichiers, puis exécutez-les dans un *shell*.

7.6.1 Lecture et saut de ligne

Dans l'exemple 'girafe', 'tigre', etc, ci-dessus, comment expliquez-vous que Python saute une ligne à chaque itération ? Réécrivez les instructions *ad-hoc* pour que Python écrive le contenu du fichier sans sauter de ligne.

7.6.2 Écriture et saut de ligne

En reprenant le dernier exemple sur l'écriture dans un fichier, vous pouvez constater que les noms d'animaux ont été écrits les uns à la suite des autres, sans retour à la ligne. Comment expliquez-vous ce résultat ? Modifiez les instructions de manière à écrire un animal par ligne.

7.6.3 Structures secondaires

Dans cet exercice, nous allons utiliser une sortie partielle de DSSP (*Define Secondary Structure of Proteins*), qui est un logiciel d'assignation des structures secondaires des protéines. Ce fichier contient 5 colonnes correspondant respectivement au numéro de résidu, à l'acide aminé, sa structure secondaire et ses angles phi/psi.

- Téléchargez le fichier `first_helix_1tfe.txt`³ sur le site de notre cours et sauvegardez-le dans votre répertoire de travail (jetez-y un oeil en passant).
- Chargez les lignes de ce fichier en les plaçant dans une liste puis fermez le fichier.
- Écrivez chaque ligne à l'écran pour vérifier que vous avez bien chargé le fichier.
- Écrivez dans un fichier `output.txt` chacune des lignes. N'oubliez pas le retour à la ligne pour chaque acide aminé.
- Écrivez dans un fichier `output2.txt` chacune des lignes suivies du message `line checked` suivi d'un saut à la ligne.

7.6.4 Spirale (exercice +++)

On se propose de réaliser un script qui écrit dans un fichier `spirale.dat` les coordonnées cartésiennes d'une spirale. *Conseil* : si vous faites varier une variable `theta` qui représente un angle entre 0 à 2π

3. https://python.sdv.univ-paris-diderot.fr/data-files/first_helix_1tfe.txt

(ou un multiple de 2π), les fonctions cosinus et sinus vous permettront d'obtenir les coordonnées des points sur un cercle.

Si vous avez réussi cette première étape, tentez la spirale en faisant en sorte que le rayon du cercle augmente au fur et à mesure que `theta` grandit.

Pour admirer votre spirale, vous pourrez utiliser le programme `xmgrace` en tapant

```
1 xmgrace spirale.dat
```

dans un shell bash.

Chapitre 8

Modules

8.1 Définition

Les modules sont des programmes Python qui contiennent des fonctions que l'on est amené à réutiliser souvent (on les appelle aussi bibliothèques ou *libraries*). Les développeurs de Python ont mis au point de nombreux modules qui effectuent une quantité phénoménale de tâches. Pour cette raison, prenez le réflexe de vérifier si une partie du code que vous souhaitez écrire n'existe pas déjà sous forme de module. La plupart de ces modules sont déjà installés dans les versions standards de Python. Vous pouvez accéder à une documentation exhaustive¹ sur le site de Python. Explorez un peu ce site, la quantité de modules disponibles est impressionnante (plus de 300).

8.2 Importation de modules

Jusqu'à présent, nous avons rencontré une fois cette notion de module lorsque nous avons voulu tirer un nombre aléatoire.

```
1 >>> import random
2 >>> random.randint(0,10)
3 4
```

Regardons de plus près cet exemple :

- L'instruction `import` permet d'accéder à toutes les fonctions du module `random`².
- Ensuite, nous utilisons la fonction (ou méthode) `randint(a,b)`³ du module `random`. Attention cette fonction renvoie un nombre entier aléatoirement tiré entre `a` inclus et `b` inclus (contrairement à la fonction `range()` par exemple). Remarquez la notation objet `random.randint()` où la fonction `randint()` peut être considérée comme une méthode de l'objet `random`.

Il existe un autre moyen d'importer une ou des fonctions d'un module :

```
1 >>> from random import randint
2 >>> randint(0,10)
3 7
```

À l'aide du mot-clé `from`, vous pouvez importer une fonction spécifique d'un module donné. Remarquez que dans ce cas il est inutile de répéter le nom du module, seul le nom de ladite fonction est requis.

1. <https://docs.python.org/3/py-modindex.html>
2. <https://docs.python.org/3/library/random.html#module-random>
3. <https://docs.python.org/3/library/random.html#random.randint>

On peut également importer toutes les fonctions d'un module :

```

1 >>> from random import *
2 >>> x = [1, 2, 3, 4]
3 >>> shuffle(x)
4 >>> x
5 [2, 3, 1, 4]
6 >>> shuffle(x)
7 >>> x
8 [4, 2, 1, 3]
9 >>> randint(0,50)
10 46
11 >>> uniform(0,2.5)
12 0.64943174760727951

```

Comme vous l'avez deviné, l'instruction `from random import *` importe toutes les fonctions du module `random`. On peut ainsi utiliser toutes ses fonctions directement, comme par exemple `shuffle()`, qui permute une liste aléatoirement.

Dans la pratique, plutôt que de charger toutes les fonctions d'un module en une seule fois (`from random import *`), nous vous conseillons de charger le module (`import random`) puis d'appeler explicitement les fonctions voulues (`random.randint(0,2)`).

Il est également possible de définir un alias (un nom plus court) pour un module :

```

1 >>> import random as rand
2 >>> rand.randint(1, 10)
3 6
4 >>> rand.uniform(1, 3)
5 2.643472616544236

```

Dans cet exemple, les fonctions du module `random` sont accessibles via l'alias `rand`.

Enfin, si vous voulez vider de la mémoire un module déjà chargé, vous pouvez utiliser l'instruction `del` :

```

1 >>> import random
2 >>> random.randint(0,10)
3 2
4 >>> del random
5 >>> random.randint(0,10)
6 Traceback (most recent call last):
7   File "<stdin>", line 1, in ?
8 NameError: name 'random' is not defined

```

Vous constatez qu'un rappel d'une fonction du module `random` après l'avoir vidé de la mémoire retourne un message d'erreur.

8.3 Obtenir de l'aide sur les modules importés

Pour obtenir de l'aide sur un module rien de plus simple, il suffit d'invoquer la commande `help()` :

```

1 >>> import random
2 >>> help(random)
3 ...

```

On peut se déplacer dans l'aide avec les flèches ou les touches *page-up* et *page-down* (comme dans les commandes Unix *man*, *more* ou *less*). Il est aussi possible d'invoquer de l'aide sur une fonction particulière d'un module de la manière suivante `help(random.randint)`.

La commande `help()` est en fait une commande plus générale permettant d'avoir de l'aide sur n'importe quel objet chargé en mémoire.

```

1 >>> t = [1, 2, 3]
2 >>> help(t)
3 Help on list object:
4
5 class list(object)
6 | list() -> new list
7 | list(sequence) -> new list initialized from sequence's items
8 |
9 | Methods defined here:
10 |
11 | __add__(...)
12 |     x.__add__(y) <==> x+y
13 |
14 ...

```

Si on veut connaître d'un seul coup d'oeil toutes les méthodes ou variables associées à un objet, on peut utiliser la fonction `dir()` :

```

1 >>> import random
2 >>> dir(random)
3 ['BPF', 'LOG4', 'NV_MAGICCONST', 'RECIP_BPF', 'Random', 'SG_MAGICCONST', 'SystemRandom',
4 'TWOPI', 'WichmannHill', '_BuiltinMethodType', '_MethodType', '__all__', '__builtins__',
5 '__doc__', '__file__', '__name__', '_acos', '_ceil', '_cos', '_e', '_exp', '_hexlify',
6 '_inst', '_log', '_pi', '_random', '_sin', '_sqrt', '_test', '_test_generator',
7 '_urandom', '_warn', 'betavariate', 'choice', 'expovariate', 'gammavariate', 'gauss',
8 'getrandbits', 'getstate', 'jumpahead', 'lognormvariate', 'normalvariate',
9 'paretovariate', 'randint', 'random', 'randrange', 'sample', 'seed', 'setstate',
10 'shuffle', 'uniform', 'vonmisesvariate', 'weibullvariate']
11 >>>

```

8.4 Modules courants

Il existe une série de modules que vous serez probablement amenés à utiliser si vous programmez en Python. En voici une liste non exhaustive. Pour la liste complète, reportez-vous à la page des modules⁴ sur le site de Python :

- `math`⁵ : fonctions et constantes mathématiques de base (sin, cos, exp, pi...).
- `sys`⁶ : passage d'arguments, gestion de l'entrée/sortie standard...
- `os`⁷ : dialogue avec le système d'exploitation (permet de sortir de Python, lancer une commande en *shell*...).
- `random`⁸ : génération de nombres aléatoires.

4. <https://docs.python.org/3/py-modindex.html>

5. <https://docs.python.org/3/library/math.html#module-math>

6. <https://docs.python.org/3/library/sys.html#module-sys>

7. <https://docs.python.org/3/library/os.html#module-os>

8. <https://docs.python.org/3/library/random.html#module-random>

- `time`⁹ : permet d'accéder à l'heure de l'ordinateur et aux fonctions gérant le temps.
- `calendar`¹⁰ : fonctions de calendrier.
- `profile`¹¹ : permet d'évaluer le temps d'exécution de chaque fonction dans un programme (*profiling* en anglais).
- `urllib`¹² : permet de récupérer des données sur internet depuis Python.
- `tkinter`¹³ : interface python avec Tk (permet de créer des objets graphiques)
- `re`¹⁴ : gestion des expressions régulières.
- `pickle`¹⁵ : écriture et lecture de structures Python (comme les dictionnaires par exemple).

Nous vous conseillons vivement d'aller surfer sur les pages de ces modules pour découvrir toutes leurs potentialités.

Vous verrez plus tard comment créer votre propres modules lorsque vous êtes amenés à réutiliser souvent vos propres fonctions.

Enfin, notez qu'il existe de nombreux autres modules qui ne sont pas installés de base dans Python mais qui sont de grand intérêt en bioinformatique (au sens large). Citons-en quelques-uns : `numpy` (notion de matrice, algèbre linéaire, transformée de Fourier), `biopython` (recherche dans les banques de données biologiques), `matplotlib` (construction de graphiques)...

8.5 Module sys : passage d'arguments

Le `sys`¹⁶ contient (comme son nom l'indique) des fonctions et des variables spécifiques au système, ou plus exactement à l'interpréteur lui-même. Par exemple, il permet de gérer l'entrée (*stdin*) et la sortie standard (*stdout*). Ce module est particulièrement intéressant pour récupérer les arguments passés à un script Python lorsque celui-ci est appelé en ligne de commande. Dans cet exemple, écrivons le court script suivant que l'on enregistrera sous le nom `test.py` (n'oubliez pas de le rendre exécutable) :

```
1 #!/usr/bin/env python3
2
3 import sys
4 print(sys.argv)
```

Ensuite lancez `test.py` suivi de plusieurs arguments. Par exemple :

```
1 $ python3 test.py salut girafe 42
2 ['test.py', 'salut', 'girafe', '42']
```

Dans l'exemple précédent, `$` représente l'invite du *shell* Linux, `test.py` est le nom du script Python, `salut`, `girafe` et `42` sont les arguments passés au script.

La variable `sys.argv` est une liste qui représente tous les arguments de la ligne de commande, y compris le nom du script lui même qu'on peut retrouver dans `sys.argv[0]`. On peut donc accéder à chacun de ces arguments avec `sys.argv[1]`, `sys.argv[2]`...

On peut aussi utiliser la fonction `sys.exit()` pour quitter un script Python. On peut donner un argument à cette fonction (en général une chaîne de caractères) qui sera renvoyé au moment où Python

9. <https://docs.python.org/3/library/time.html#module-time>
10. <https://docs.python.org/3/library/calendar.html#module-calendar>
11. <https://docs.python.org/3/library/profile.html#module-profile>
12. <https://docs.python.org/3/library/urllib.html#module-urllib>
13. <https://docs.python.org/3/library/tkinter.html#module-tkinter>
14. <https://docs.python.org/3/library/re.html#module-re>
15. <https://docs.python.org/3/library/pickle.html#module-pickle>
16. <https://docs.python.org/3/library/sys.html#module-sys>

quittera le script. Par exemple, si vous attendez au moins un argument en ligne de commande, vous pouvez renvoyer un message pour indiquer à l'utilisateur ce que le script attend comme argument :

```

1 #!/usr/bin/env python3
2
3 import sys
4
5 if len(sys.argv) != 2:
6     sys.exit("ERREUR : il faut exactement un argument.")
7
8 #
9 # suite du script
10 #

```

Puis on l'exécute sans argument :

```

1 $ python3 test.py
2 ERREUR : il faut exactement un argument.

```

Notez qu'ici on vérifie que le script possède deux arguments car le nom du script lui-même est le premier argument et `file.txt` constitue le second.

8.6 Module os

Le module `os` gère l'interface avec le système d'exploitation.

`os.path.exists()` est une fonction pratique de ce module qui vérifie la présence d'un fichier sur le disque.

```

1 >>> import sys
2 >>> import os
3 >>> if os.path.exists("toto.pdb"):
4 ...     print("le fichier est présent")
5 ... else:
6 ...     sys.exit("le fichier est absent")
7 ...
8 le fichier est absent

```

Dans cet exemple, si le fichier n'est pas présent sur le disque, on quitte le programme avec la fonction `exit()` du module `sys`.

La fonction `system()` permet d'appeler n'importe quelle commande externe.

```

1 >>> import os
2 >>> os.system("ls -al")
3 total 5416
4 drwxr-xr-x 2 poulain bioinfo 4096 2010-07-21 14:33 .
5 drwxr-xr-x 6 poulain bioinfo 4096 2010-07-21 14:26 ..
6 -rw-r--r-- 1 poulain bioinfo 124335 2010-07-21 14:31 1BTA.pdb
7 -rw-r--r-- 1 poulain bioinfo 4706057 2010-07-21 14:31 NC_000913.fna
8 -rw-r--r-- 1 poulain bioinfo 233585 2010-07-21 14:30 NC_001133.fna
9 -rw-r--r-- 1 poulain bioinfo 463559 2010-07-21 14:33 NC_001133.gbk
10 0

```

La commande externe `ls -al` est introduite comme une chaîne de caractères à la fonction `system()`. La valeur 0 sur la dernière ligne correspondant au code de sortie de la commande `ls -al` et indique ici qu'il n'y a pas eu d'erreur. Attention à ne pas faire de bêtise avec cette fonction qui peut exécuter n'importe quelle commande Unix.

8.7 Exercices

Conseil : pour les trois premiers exercices, utilisez l'interpréteur Python. Pour les exercices suivants, écrivez des scripts dans des fichiers, puis exécutez-les dans un *shell*.

8.7.1 Racine carrée

Affichez sur la même ligne les nombres de 10 à 20 (inclus) ainsi que leur racine carrée avec 3 décimales. Utilisez pour cela le module `math` avec la fonction `sqrt()` (documentation¹⁷). Exemple :

```
1 10 3.162
2 11 3.317
3 12 3.464
4 13 3.606
5 ...
```

8.7.2 Cosinus

Calculez le cosinus de $\pi/2$ en utilisant le module `math` avec la fonction `cos()` (documentation¹⁸) et la constante `pi` (documentation¹⁹).

8.7.3 Liste de fichiers dans un répertoire

Affichez la liste des fichiers du répertoire courant avec le module `os`. N'utilisez pas la fonction `os.system()` mais la fonction `os.listdir()`. Lisez la documentation²⁰ pour comprendre comment l'utiliser.

8.7.4 Affichage temporisé

Affichez les nombres de 1 à 10 avec 1 seconde d'intervalle. Utilisez pour cela le module `time` et sa fonction `sleep()` (documentation²¹).

8.7.5 Séquences aléatoires de chiffres

Générez une séquence aléatoire de 6 chiffres, ceux-ci étant des entiers tirés entre 1 et 4. Utilisez le module `random` et la fonction `randint()` (documentation²²).

17. <https://docs.python.org/3/library/math.html#math.sqrt>

18. <https://docs.python.org/3/library/math.html#math.cos>

19. <https://docs.python.org/3/library/math.html#math.pi>

20. <https://docs.python.org/3/library/os.html#os.listdir>

21. <https://docs.python.org/3/library/time.html#time.sleep>

22. <https://docs.python.org/3/library/random.html#random.randint>

8.7.6 Séquences aléatoires de bases

Générez une séquence aléatoire de 20 bases de deux manières différentes. Utilisez le module `random` avec la fonction `randint()` (documentation [23](#)) ou `choice()` (documentation [24](#)).

8.7.7 Jour de naissance

Déterminez votre jour (lundi, mardi...) de naissance en utilisant la fonction `weekday()` (documentation [25](#)) du module `calendar`.

8.7.8 Détermination du nombre pi par la méthode Monte Carlo (exercice +++)

Soit un cercle de rayon 1 (en trait plein sur la figure 8.1) inscrit dans un carré de côté 2 (en trait pointillé).

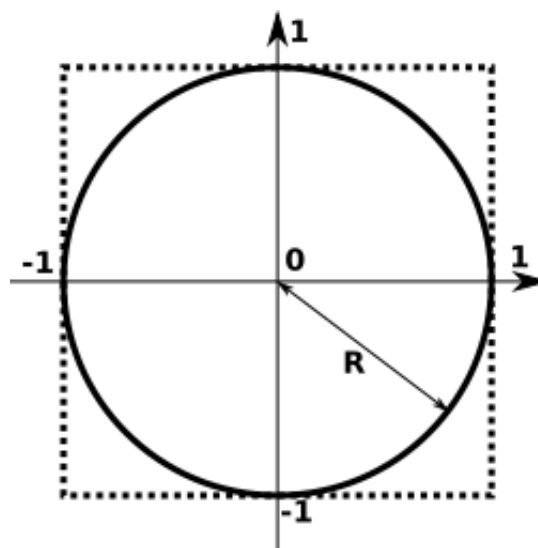


FIGURE 8.1 – Cercle de rayon 1 inscrit dans un carré de côté 2.

Avec $R = 1$, l'aire du carré vaut $(2R)^2$ soit 4 et l'aire du cercle vaut πR^2 soit π .

En choisissant N points aléatoires (à l'aide d'une distribution uniforme) à l'intérieur du carré, la probabilité que ces points se trouvent aussi dans le cercle est

$$p = \frac{\text{aire du cercle}}{\text{aire du carré}} = \frac{\pi}{4}$$

Soit n , le nombre points effectivement dans le cercle, il vient alors

$$p = \frac{n}{N} = \frac{\pi}{4},$$

d'où

$$\pi = 4 \times \frac{n}{N}.$$

23. <https://docs.python.org/3/library/random.html#random.randint>

24. <https://docs.python.org/3/library/random.html#random.choice>

25. <https://docs.python.org/3/library/calendar.html#calendar.weekday>

Déterminez une approximation de π par cette méthode. Pour cela, pour N itérations :

- Choisissez aléatoirement les coordonnées x et y d'un point entre -1 et 1. Utilisez la fonction `uniform()` (documentation²⁶) du module `random`.
- Calculez la distance entre le centre du cercle et ce point.
- Déterminez si cette distance est inférieure au rayon du cercle, c'est-à-dire si le point est dans le cercle ou pas.
- Si le point est effectivement dans le cercle, incrémentez le compteur n .

Finalement calculez le rapport entre n et N et proposer une estimation de π . Quelle valeur de π obtenez-vous pour 100 itérations? 1000 itérations? 10000 itérations? Comparez les valeurs obtenues à la valeur de π fournie par le module `math`.

On rappelle que la distance d entre deux points A et B de coordonnées respectives (x_A, y_A) et (x_B, y_B) se calcule comme :

$$d = \sqrt{(x_B - x_A)^2 + (y_B - y_A)^2}$$

26. <https://docs.python.org/3/library/random.html#random.uniform>

Chapitre 9

Fonctions

9.1 Principe et généralités

En programmation, les **fonctions** sont très utiles pour réaliser plusieurs fois la même opération au sein d'un programme. Elles permettent également de rendre le code plus lisible et plus clair en le fractionnant en blocs logiques.

Vous connaissez déjà certaines fonctions Python, par exemple `math.cos(angle)` du module `math` renvoie le cosinus de la variable `angle` exprimé en radian. Vous connaissez aussi des fonctions internes à Python comme `range()` ou `len()`. Pour l'instant, une fonction est à vos yeux une sorte de *boîte noire* (cf Figure 9.1) :

- À laquelle vous passez une (ou zero ou plusieurs) valeur(s) entre parenthèses. Ces valeurs sont appelées arguments.
- Qui effectue une action. Par exemple la fonction `random.shuffle()` permute aléatoirement une liste.
- Et qui renvoie éventuellement un résultat (plus précisément un objet).

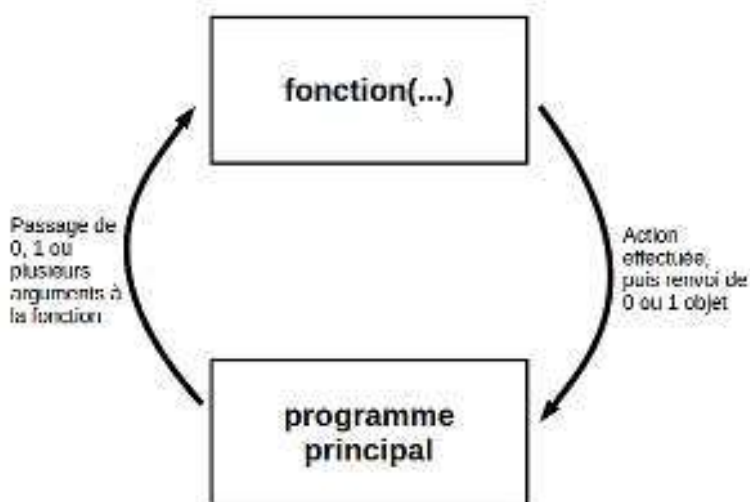


FIGURE 9.1 – Fonctionnement schématique d'une fonction

Par exemple si vous appelez la fonction `len()` en lui passant l'argument `[0, 1, 2]` :

```
1 >>> len([0, 1, 2])
2 3
```

celle-ci vous renvoie un entier indiquant la longueur de la liste passée en argument.

Autre exemple, si vous appelez la méthode `liste.append(5)` (n'oubliez pas, une méthode est une **fonction** qui agit sur l'objet auquel elle est attachée), `append()` ajoute l'entier 5 à la variable `liste` mais ne renvoie rien.

Au contraire, aux yeux du programmeur une fonction est une portion de code effectuant une suite d'instructions bien particulière. Avant de démarrer sur la syntaxe, revenons sur cette notion de *boîte noire* :

- Une fonction effectue une tâche. Pour cela, elle reçoit éventuellement des arguments et renvoie éventuellement un résultat. L'algorithme utilisé au sein de la fonction se n'intéresse pas directement l'utilisateur. Par exemple, il est inutile de savoir comment la fonction `math.cos()` calcule un cosinus, on a juste besoin de savoir qu'il faut lui passer en argument un angle en radian et qu'elle renvoie le cosinus de cet angle. Ce qu'il se passe à l'intérieur de la fonction ne regarde que le programmeur (c'est-à-dire vous dans ce chapitre).
- Chaque fonction effectue en général une tâche **unique et précise**. Si cela se complique, il est plus judicieux d'écrire plusieurs fonctions (qui peuvent éventuellement s'appeler les unes les autres). Cette **modularité** améliore la qualité générale et la lisibilité du code. Vous verrez qu'en Python, les fonctions présentent une grande flexibilité.

Pour finir sur les généralités, vous avez vu dans la Figure 9.1 que nous avons utilisé le terme **programme principal** (*main* en anglais) pour désigner l'endroit depuis lequel on appelle une fonction (on verra plus tard que l'on peut en fait appeler une fonction de n'importe où). Le programme principal désigne le code qui est exécuté lorsqu'on lance le script Python, c'est à dire toute la suite d'instructions qui commencent à la colonne 1 du script, autrement dit toutes les instructions en dehors des fonctions. En général, dans un script Python, on écrit d'abord les fonctions puis le programme principal (donc celui-ci se situe à la fin du script). Nous aurons l'occasion de revenir sur cette notion de programme principal dans la section 9.4 de ce chapitre ainsi que dans le chapitre 12.

9.2 Définition

Pour définir une fonction, Python utilise le mot-clé `def` et si on veut que celle-ci renvoie une valeur, il faut utiliser le mot-clé `return`. Par exemple :

```

1 >>> def carre(x):
2 ...     return x**2
3 ...
4 >>> print(carre(2))
5 4

```

Notez que la syntaxe de `def` utilise les `:` comme les boucles `for` et `while` ainsi que les tests `if`, un bloc d'instructions est donc attendu. De même que pour les boucles et les tests, l'**indentation** de ce bloc d'instructions (*i.e.* le corps de la fonction) est **obligatoire**.

Dans l'exemple précédent, nous avons passé un argument à la fonction `carre()` qui nous a retourné une valeur que nous avons affichée à l'écran. Que veut dire valeur retournée ? Et bien cela signifie que cette dernière est stockable dans une variable :

```

1 >>> res = carre(2)
2 >>> print(res)
3 4

```

Ici, le résultat renvoyé par la fonction est stockée dans la variable `res`. Notez qu'une fonction ne prend pas forcément un argument et ne renvoie pas forcément une valeur, par exemple :

```

1 >>> def hello():
2 ...     print("bonjour")
3 ...
4 >>> hello()
5 bonjour

```

Dans ce cas la fonction `hello()` se contente d'imprimer la chaîne de caractères `"hello"` à l'écran. Elle ne prend aucun argument et ne renvoie aucun résultat. Par conséquent, cela n'a pas de sens de vouloir récupérer dans une variable le résultat renvoyé par une telle fonction. Si on essaie tout de même, Python affecte la valeur `None` qui signifie *rien* en anglais :

```

1 >>> x = hello()
2 bonjour
3 >>> print(x)
4 None

```

Ceci n'est pas une faute car Python n'émet pas d'erreur, toutefois cela ne présente, la plupart du temps, guère d'intérêt.

9.3 Passage d'arguments

Le nombre d'argument(s) que l'on peut passer à une fonction est variable. Nous avons vu ci-dessus des fonctions auxquelles on passait 0 ou 1 argument. Dans les chapitres précédentes, vous avez vu des fonctions internes à Python qui prenaient au moins 2 arguments, pour rappel souvenez-vous de `range(1,10)` ou encore `range(1,10,2)`. Le nombre d'argument est donc laissé libre à l'initiative du programmeur qui est en train de développer une nouvelle fonction.

Une particularité des fonctions en Python est que vous n'êtes pas obligé de préciser le type des arguments que vous lui passez, dès lors que les opérations que vous effectuez avec ces arguments sont valides. Python est en effet connu comme étant un langage au *typage dynamique*, c'est-à-dire qu'il reconnaît pour vous le type des variables au moment de l'exécution, par exemple :

```

1 >>> def fois(x,y):
2 ...     return x*y
3 ...
4 >>> fois(2,3)
5 6
6 >>> fois(3.1415,5.23)
7 16.430045000000003
8 >>> fois('to',2)
9 'toto'

```

L'opérateur `*` reconnaît plusieurs types (entiers, réels, chaînes de caractères), notre fonction est donc capable d'effectuer des tâches différentes! Même si Python permet cela, méfiez-vous tout de même de cette grande flexibilité qui pourrait mener à des surprises dans vos futurs programmes. En général il est plus judicieux que chaque argument ait un type précis (*int*, *str*, *float*, etc), et pas l'un ou l'autre.

Un énorme avantage en Python est que les fonctions sont capables de renvoyer plusieurs valeurs à la fois, comme dans cette fraction de code :

```

1 >>> def carre_cube(x):
2 ...     return x**2,x**3
3 ...

```

```

4 >>> carre_cube(2)
5 (4, 8)

```

En réalité Python ne renvoie qu'un seul objet, mais celui-ci peut être séquentiel, c'est à dire contenir lui même plusieurs objets. Dans notre exemple Python renvoie un objet de type `tuple`, type que nous verrons dans le chapitre 13 (il s'agit d'une sorte de liste avec des propriétés différentes). Notre fonction pourrait tout autant renvoyer une liste :

```

1 >>> def carre_cube2(x):
2 ...     return [x**2,x**3]
3 ...
4 >>> carre_cube2(3)
5 [9, 27]

```

Renvoyer un *tuple* ou une liste de deux arguments (ou plus) est notamment très pratique en conjonction avec l'**affectation multiple**, par exemple :

```

1 >>> z1,z2 = carre_cube2(3)
2 >>> z1
3 9
4 >>> z2
5 27

```

Cela permet de récupérer plusieurs valeurs retournées par une fonction et les affecter à des variables différentes à la volée.

Enfin, il est possible de passer un ou plusieurs argument(s) de manière facultative et de leur attribuer une valeur par défaut :

```

1 >>> def useless_fct(x=1):
2 ...     return x
3 ...
4 >>> useless_fct()
5 1
6 >>> useless_fct(10)
7 10

```

Notez que si on passe plusieurs arguments à une fonction, le ou les arguments facultatifs doivent être situés après les arguments obligatoires. Il faut donc écrire `def fct(x, y, z=1):`.

9.4 Variables locales et variables globales

Lorsqu'on manipule les fonctions il est essentiel de bien comprendre comment se comportent les variables. Une variable est dite *locale* lorsqu'elle est créée dans une fonction, car elle n'existera et ne sera visible que lors de l'exécution de la dite fonction. Une variable est dite *globale* lorsqu'elle est créée dans le programme principal; elle sera visible partout dans le programme.

Ceci ne vous paraît pas clair? Nous allons suivre un exemple simple illustrant notre propos et qui vous permettra de saisir aisément ces concepts. Regardez le code suivant :

```

1 # définition d'une fonction carre()
2 def carre(x):
3     y = x**2
4     return y

```

```

5
6 # programme principal
7 z = 5
8 resultat = carre(5)
9 print(resultat)

```

Pour la suite nous allons utilisé l'excellent site [pythontutor](http://www.pythontutor.com)¹ qui permet de visualiser l'état des variables au fur et à mesure de l'exécution d'un code Python (avant de continuer nous vous conseillons de prendre 5 minutes pour tester ce site). Regardons maintenant ce qu'il se passe dans le code ci-dessus, étape par étape :

- Étape 1 : Python est prêt à lire la première ligne de code.

Python 3.6

```

→ 1 def carre(x):
2     y = x**2
3     return y
4
5 # programme principal
6 z = 5
7 resultat = carre(5)
8 print(resultat)

```

Print output (drag lower right corner to resize)

Frames Objects

- Étape 2 : Python met en mémoire la fonction `carre()` (notez qu'il ne l'exécute pas!). La fonction est mise dans une case de la mémoire nommée *global frame* qui est l'espace mémoire du programme principal. Y seront stockées toutes les variables *globales* créées dans le programme. Python est maintenant prêt à exécuter le programme principal.

Python 3.6

```

→ 1 def carre(x):
2     y = x**2
3     return y
4
5 # programme principal
→ 6 z = 5
7 resultat = carre(5)
8 print(resultat)

```

Print output (drag lower right corner to resize)

Frames Objects

Global frame

carre → function carre(x)

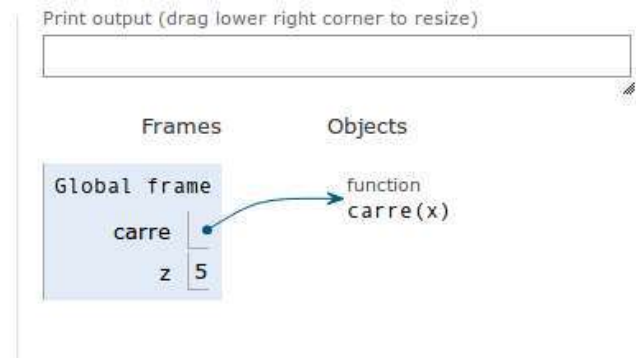
- Étape 3 : Python lit et met en mémoire la variable `z`. Celle-ci étant créée dans le programme principal, il s'agira d'une variable *globale*. Ainsi, elle sera également stockée dans le *global frame*.

1. <http://www.pythontutor.com>

```

Python 3.6
1 def carre(x):
2     y = x**2
3     return y
4
5 # programme principal
6 z = 5
7 resultat = carre(5)
8 print(resultat)

```

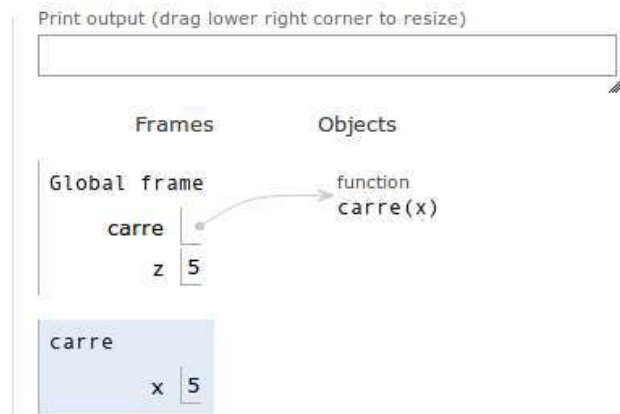


- Étape 4 : La fonction `carre()` est appelée et on lui passe en argument l'entier 5. La fonction rentre alors en exécution et un nouveau cadre bleu est créé dans lequel *pythontutor* va nous indiquer toutes les variables *locales* à la fonction. Notez bien que la variable passée en argument, qui s'appelle `x` dans la fonction, est créée en tant que variable *locale*. On pourra aussi remarquer que les variables *globales* situées dans le *global frame* sont toujours là.

```

Python 3.6
1 def carre(x):
2     y = x**2
3     return y
4
5 # programme principal
6 z = 5
7 resultat = carre(5)
8 print(resultat)

```

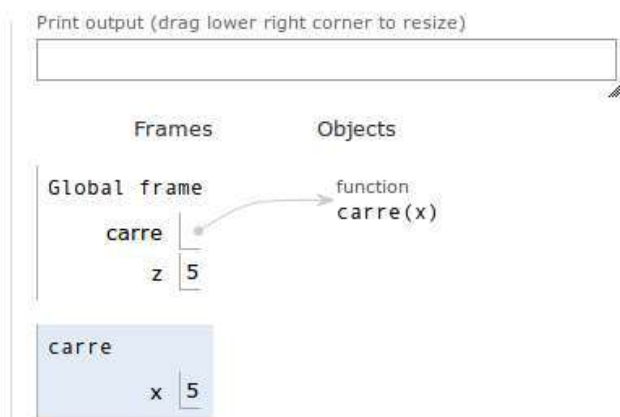


- Étape 5 : Python est maintenant prêt à exécuter chaque ligne de code de la fonction.

```

Python 3.6
1 def carre(x):
2     y = x**2
3     return y
4
5 # programme principal
6 z = 5
7 resultat = carre(5)
8 print(resultat)

```



- Étape 6 : La variable `y` est créée dans la fonction. Celle-ci est donc stockée en tant que variable *locale* à la fonction.

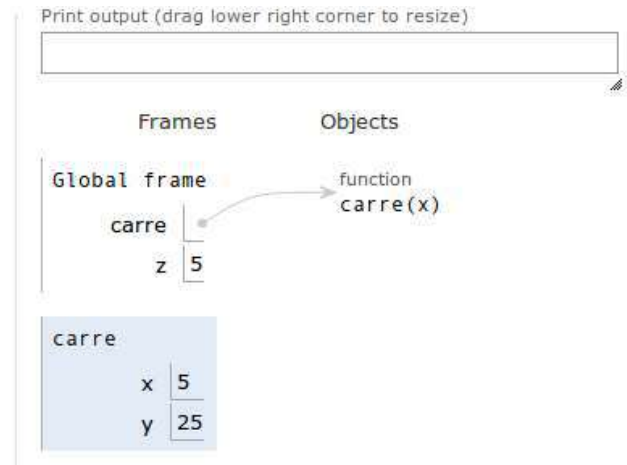
```

Python 3.6
1 def carre(x):
2     y = x**2
3     return y
4
5 # programme principal
6 z = 5
7 resultat = carre(5)
8 print(resultat)

```

[Edit code](#) | [Live programming](#)

ecuted



- Étape 7 : Python s'apprête à retourner la variable *locale* `y` au programme principal (*pythontutor* nous indique le contenu de la valeur retournée).

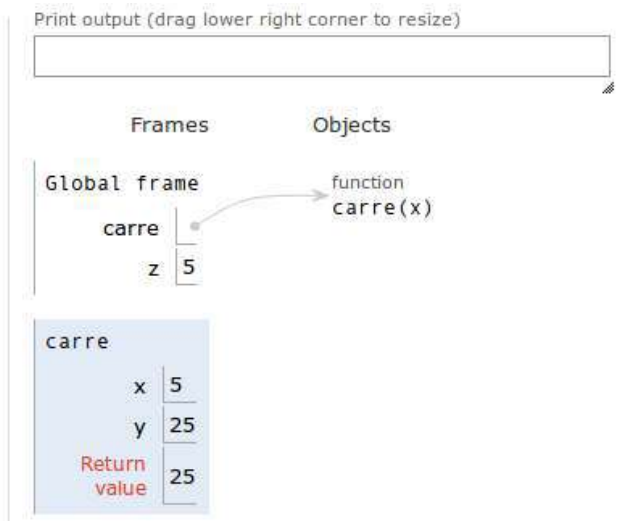
```

Python 3.6
1 def carre(x):
2     y = x**2
3     return y
4
5 # programme principal
6 z = 5
7 resultat = carre(5)
8 print(resultat)

```

[Edit code](#) | [Live programming](#)

ecuted



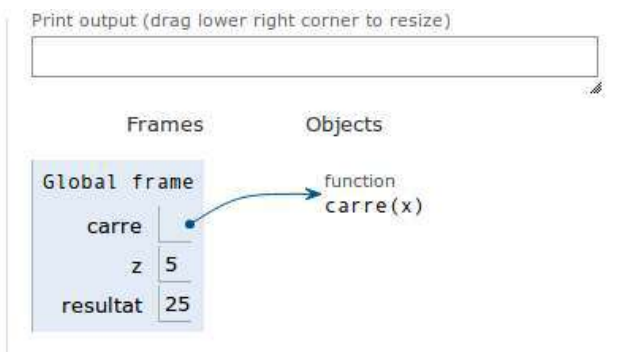
a breakpoint; use the Back and Forward buttons to jump there.

- Étape 8 : Python quitte la fonction et la valeur retournée par celle-ci est affectée à la variable *globale* `resultat`. Notez bien que lorsque Python quitte la fonction, **l'espace des variables allouées à la fonction est détruit**. Ainsi toutes les variables créées dans la fonction n'existent plus. On comprend pourquoi elles portent le nom de *locales* puisqu'elles n'existent que lorsque la fonction est en exécution.

```

Python 3.6
1 def carre(x):
2     y = x**2
3     return y
4
5 # programme principal
6 z = 5
7 resultat = carre(5)
8 print(resultat)

```

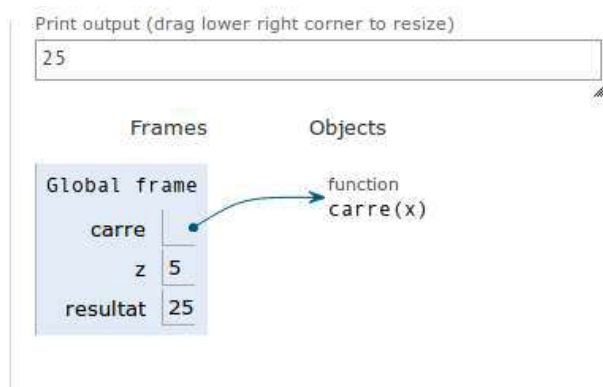


- Étape 9 : Python affiche le contenu de la variable `resultat` et l'exécution est terminée.

```

Python 3.6
1 # definition d'une fonction
2 def carre(x):
3     y = x**2
4     return y
5
6 # programme principal
7 z = 5
8 resultat = carre(5)
9 print(resultat)

```



Nous espérons que cet exemple guidé facilitera la compréhension des concepts de variables locales et globales. Cela viendra aussi avec la pratique. Nous irons un peu plus loin sur les fonctions dans le chapitre 12. D'ici là essayez de vous entraîner au maximum avec les fonctions. C'est un concept ardu, mais il est impératif de savoir les maîtriser.

Enfin, comme vous avez vu, *pythontutor* nous a grandement aidés à comprendre ce qu'il se passait. N'hésitez pas à l'utiliser sur des exemples ponctuels, celui-ci pourra vous aider à comprendre lorsqu'un code ne fait pas ce que vous attendez.

9.5 Exercices

Conseil : pour ces exercices, écrivez des scripts dans des fichiers, puis exécutez-les dans un *shell*.

9.5.1 Fonctions et pythontutor

Reprenez l'exemple de la section 9.4 à l'aide du site [pythontutor](http://www.pythontutor.com)². Regardez ensuite le code suivant et tentez de prédire sa sortie :

```

1 def calc_factorielle(n):
2     fact = 1
3     for i in range(2,n+1):
4         fact = fact * i
5     return fact
6
7 # prog principal
8 nb = 4
9 factorielle_nb = calc_factorielle(nb)
10 print("{}! = {}".format(nb,factorielle_nb))
11 nb2 = 10
12 print("{}! = {}".format(nb2,calc_factorielle(nb2)))

```

Testez ensuite cette portion de code avec *pythontutor* en cherchant à bien comprendre chaque étape. Avez-vous réussi à prédire la sortie correctement ?

2. <http://www.pythontutor.com>

9.5.2 Fonction puissance

La fonction `math.pow(x,y)` permet de calculer x^y . Reprogrammer une fonction `calc_puissance(x,y)` qui renvoie x^y en utilisant l'opérateur `**`. Dans le programme principal, calculez et affichez à l'écran 2^i avec i variant de 0 à 20. On veut que le résultat soit présenté avec le formatage suivant :

```

1 2^ 0 =      1
2 2^ 1 =      2
3 2^ 2 =      4
4 [...]
5 2^20 = 1048576

```

9.5.3 Fonction pyramide

Reprenez l'exercice 5.4.10 qui dessine un triangle isocèle. Concevez une fonction `gen_pyramide()` à laquelle vous passez un nombre entier N et qui renvoie une pyramide de N lignes (i.e. triangle isocèle) sous forme de chaîne de caractères. Le programme principal demandera à l'utilisateur le nombre de lignes souhaitées (pensez à la fonction `input()` !) et affichera la pyramide à l'écran.

9.5.4 Fonction nombre premier

Reprenez l'exercice 6.6.10 sur les nombres premiers. Concevez une fonction `is_prime()` qui prend en argument un nombre entier positif n (supérieur à 2) et qui renvoie un booléen `True` si n est premier, et `False` si n n'est pas premier. Déterminez tous les nombres premiers de 2 à 100. On souhaite avoir une sortie comme ça :

```

1  2 est premier
2  3 est premier
3  4 n'est pas premier
4  [...]
5  10 n'est pas premier
6  [...]
7  100 n'est pas premier

```

9.5.5 Fonction complement

Reprenez l'exercice 6.6.2. Créez une fonction `complement()` qui prend une liste de bases et qui renvoie la séquence complémentaire d'une séquence d'ADN sous forme de liste.

Dans le programme principal, à partir d'une séquence d'ADN `seq = ['A', 'T', 'C', 'G', 'A', 'T', 'C', 'G', 'A', 'T', 'C', 'G', 'C', 'T', 'G', 'C', 'T', 'A', 'G', 'C']` sur le brin principal, affichez `seq` et sa séquence complémentaire (en utilisant votre fonction `complement()`). On souhaite avoir une sortie comme suit :

```

brin direct          : 5'-XXXXXXXXX-3'
brin complémentaire: 3'-XXXXXXXXX-5'

```

Les 5' et 3' indique le sens de lecture, et les X représentent les bases A, T, G ou C.

9.5.6 Fonction distance

Créez une fonction `calc_distance3D()` qui calcule une distance euclidienne en 3 dimensions entre deux atomes. Testez votre fonction sur les 2 points A(0,0,0) et B(1,1,1). Trouvez-vous bien $\sqrt{3}$?

On rappelle que la distance euclidienne d entre deux points A et B de coordonnées cartésiennes respectives (x_A, y_A, z_A) et (x_B, y_B, z_B) se calcule comme suit :

$$d = \sqrt{(x_B - x_A)^2 + (y_B - y_A)^2 + (z_B - z_A)^2}$$

9.5.7 Fonctions distribution et stat

Créez une fonction `gen_distrib()` qui prend en argument les 3 entiers `debut`, `fin` et `n`. La fonction renverra une liste de n nombres réels aléatoires entre `debut` et `fin`. Cette liste suivra une distribution uniforme (utilisez pour cela la fonction `uniform()` du module `random`). Créez une autre fonction `calc_stat()` qui prend en argument une liste de nombres réels et qui renvoie une liste de 4 éléments contenant respectivement le minimum, le maximum, la moyenne et la médiane de la liste. Dans le programme principal, on souhaite générer 20 listes aléatoires de 100 nombres réels entre 0 et 100 et afficher les statistiques pour chacune d'entre elles (minimum, maximum, médiane et moyenne). On souhaite avoir une sortie comme suit :

```

1 Liste 1 : min = X.XXXX ; max = X.XXXX ; mediane = X.XXXX ; moyenne = X.XXXX
2 [...]
3 Liste 20 : min = X.XXXX ; max = X.XXXX ; mediane = X.XXXX ; moyenne = X.XXXX

```

Pour chaque ligne, on aura les statistiques (min, max, médiane, moyenne) associées à chacune des 20 listes. Les écarts sur les statistiques entre les différentes listes sont-ils importants ?

Relancez avec des listes de 1000 éléments, puis 10000 éléments. Les écarts changent-ils quand le nombre d'éléments par liste augmente ?

9.5.8 Fonction distance à l'origine

En reprenant votre fonction `calc_distance3D()`, faites-en une version 2D (`calc_distance2D()`). Ecrire une autre fonction `calc_dist2ori()` à laquelle vous passez en argument deux listes de nombres réels `list_x` et `list_y` représentant les coordonnées en 2D d'une fonction mathématique (par exemple x et $\sin(x)$). Cette fonction renverra une liste de nombres réels représentant la distance entre chaque point de la fonction et l'origine (de coordonnées (0,0)). Votre programme devra générer un fichier `sin2ori.dat` qui contiendra deux colonnes : la première représente les x , la deuxième la distance entre chaque point de la fonction $\sin(x)$ à l'origine. On souhaite faire ce calcul entre $-\pi$ et π par pas de 0.1. Vous pourrez admirer votre résultat avec la commande :

```
xmgrace sin2ori.dat
```

9.5.9 Fonction aire sous la courbe (exercice +++)

La méthode des trapèzes³ permet de calculer numériquement l'intégrale d'une fonction. Elle consiste à évaluer l'aire sous une fonction en évaluant l'aire de trapèzes successifs. On souhaite créer une fonction `calc_aire()` qui prend en argument deux listes de `floats` `list_x` et `list_y` représentant les coordonnées d'une fonction (par exemple x et $\sin(x)$) et qui renvoie l'aire sous la courbe. On calculera les intégrales suivantes :

3. https://fr.wikipedia.org/wiki/M%C3%A9thode_des_trap%C3%A8zes

- $\int_0^1 x dx$
- $\int_0^1 \sqrt{x} dx$
- $\int_{-\pi}^{+\pi} \sin(x) dx$
- $\int_{-\pi}^{+\pi} \cos(x) dx$

Chapitre 10

Plus sur les chaînes de caractères

10.1 Préambule

Nous avons déjà abordé les chaînes de caractères dans le chapitre *variables* et *affichage*. Ici nous allons un peu plus loin notamment avec les méthodes associées aux chaînes de caractères¹.

10.2 Chaînes de caractères et listes

Les chaînes de caractères peuvent être considérées comme des listes (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'
```

Mais *a contrario* des listes, les chaînes de caractères présentent toutefois une différence notable, ce sont **des listes non modifiables**. Une fois 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):
```

1. <https://docs.python.org/fr/3.6/library/string.html>

```

6 File "<stdin>", line 1, in <module>
7 TypeError: 'str' object does not support item assignment

```

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

10.3 Caractères spéciaux

Il existe certains caractères spéciaux comme `\n` que nous avons déjà vu (pour le retour à la ligne). Le caractère `\t` vous permet d'écrire une tabulation. Si vous voulez écrire un guillemet simple ou double (et que celui-ci ne soit pas confondu avec les guillemets de déclaration de la chaîne de caractères), vous pouvez utiliser `'` ou `"` ou utiliser respectivement des guillemets doubles ou simple pour déclarer votre chaîne de caractères.

```

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
6 >>> print("Un brin d'ADN")
7 Un brin d'ADN
8 >>> print('Python est un "super" langage de programmation')
9 Python est un "super" langage de programmation

```

Lorsqu'on souhaite écrire un texte sur plusieurs lignes, il est très commode d'utiliser les guillemets triples permettant de conserver 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

```

10.4 Méthodes associées aux chaînes de caractères

Voici quelques méthodes² spécifiques aux objets de type `string` :

```

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

```

2. <https://docs.python.org/fr/3.6/library/string.html>

```
7 'tigre'
```

Les fonctions `lower()` et `upper()` renvoient un texte en minuscule et en majuscule respectivement. On remarque que l'utilisation de ces fonctions n'altèrent pas la chaîne de départ mais renvoie la chaîne transformée.

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

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

ou encore plus simple avec la fonction Python adéquate :

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

Il existe une méthode associée aux chaînes de caractères qui est particulièrement pratique, la fonction `split()` :

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

La fonction `split()` découpe la ligne en plusieurs éléments appelés *champs*, en utilisant comme séparateur les espaces ou les tabulations. Il est possible de modifier le séparateur de champs, par exemple :

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

La fonction `find()` 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('tig')
7 -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 retourné :

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

On trouve aussi la fonction `replace()`, qui serait l'équivalent de la fonction de substitution de la commande Unix `sed` :

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

Enfin, la fonction `count()` 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
```

10.5 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 les manipuler.

Considérons par exemple la chaîne de caractères `val` :

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

Nous souhaitons extraire les valeurs 3.4 et 17.2 pour ensuite les additionner.

Dans un premier temps, nous allons découper la chaîne de caractères avec la commande `.split()` :

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

Nous obtenons une liste de chaînes de caractères. Nous devons transformer les deux premiers éléments de cette liste en réels (avec la fonction `float()`) pour pouvoir les additionner :

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

10.6 Conversion d'une liste de chaînes de caractères en une chaîne de caractères

On a vu dans le chapitre 2 la conversion des types simples (entier, réel et chaînes 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 appelle à la fonction `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, nous avons utilisé un tiret, un espace et rien.

Attention, la fonction `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

```

Nous espérons 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 >>> dir(animaux)
2 ['__add__', '__class__', '__contains__', '__delattr__', '__doc__', '__eq__', '__ge__',
3  '__getattr__', '__getitem__', '__getnewargs__', '__getslice__', '__gt__',
4  '__hash__', '__init__', '__le__', '__len__', '__lt__', '__mod__', '__mul__', '__ne__',
5  '__new__', '__reduce__', '__reduce_ex__', '__repr__', '__rmod__', '__rmul__',
6  '__setattr__', '__str__', 'capitalize', 'center', 'count', 'decode', 'encode',
7  'endswith', 'expandtabs', 'find', 'index', 'isalnum', 'isalpha', 'isdigit', 'islower',
8  'isspace', 'istitle', 'isupper', 'join', 'ljust', 'lower', 'lstrip', 'partition',
9  'replace', 'rfind', 'rindex', 'rjust', 'rpartition', 'rsplit', 'rstrip', 'split',
10 'splitlines', 'startswith', 'strip', 'swapcase', 'title', 'translate', 'upper',
11 'zfill']

```

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

Vous pouvez ensuite accéder à l'aide et à la documentation d'une fonction particulière avec `help()` :

```

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)

```

10.7 Exercices

Conseil : pour ces exercices, écrivez des scripts dans des fichiers, puis exécutez-les dans un *shell*.

10.7.1 Parcours d'une liste de chaînes de caractères

Soit la liste ['girafe', 'tigre', 'singe', 'souris']. Avec une boucle, affichez chaque élément ainsi que sa taille (nombre de caractères).

10.7.2 Fréquence des bases dans une séquence nucléique

Soit la séquence nucléique ATATACGGATCGGCTGTTGCCTGCGTAGTAGCGT. On souhaite calculer la fréquence de chaque base A, T, C et G dans cette séquence et afficher le résultat à l'écran. Créez pour cela une fonction `calc_composition()` à laquelle vous passez votre séquence nucléique sous forme de chaîne de caractères et qui renvoie une liste de 4 nombres réels (*float*) indiquant respectivement la fréquence en bases A, T, G et C.

10.7.3 Conversion des acides aminés du code à trois lettres au code à une lettre

Soit la séquence protéique ALA GLY GLU ARG TRP TYR SER GLY ALA TRP. Transformez cette séquence en une chaîne de caractères en utilisant le code à une lettre pour les acides aminés.

Rappel de la nomenclature des acides aminés :

Acide aminé	Code 3-lettres	Code 1-lettre
Alanine	Ala	A
Arginine	Arg	R
Asparagine	Asn	N
Aspartate	Asp	D
Cystéine	Cys	C
Glutamate	Glu	E
Glutamine	Gln	Q
Glycine	Gly	G
Histidine	His	H
Isoleucine	Ile	I
Leucine	Leu	L
Lysine	Lys	K
Méthionine	Met	M
Phénylalanine	Phe	F
Proline	Pro	P
Sérine	Ser	S
Thréonine	Thr	T
Tryptophane	Trp	W
Tyrosine	Tyr	Y
Valine	Val	V

10.7.4 Distance de Hamming

La distance de Hamming³ mesure la différence entre deux séquences de même taille en sommant le nombre de positions qui, pour chaque séquence, ne correspondent pas au même acide aminé.

Écrivez la fonction `hamming()` qui prend en argument deux chaînes de caractères et qui renvoie la distance de Hamming entre ces deux chaînes de caractères.

Calculez la distance de Hamming entre les séquences `AGWPSGGASAGLAIL` et `IGWPSAGASAGLWIL`, puis entre les séquences `ATTCATACGTTACGATT` et `ATACTTACGTAACCATT`.

10.7.5 Palindrome

Un palindrome est un mot ou une phrase dont l'ordre des lettres reste le même si on le lit de gauche à droite ou de droite à gauche. Par exemple, *ressasser* et *Engage le jeu que je le gagne* sont des palindromes.

Écrivez la fonction `palindrome()` qui prend en argument une chaîne de caractères et qui affiche `xxx est un palindrome` si la chaîne de caractères est un palindrome et `xxx n'est pas un palindrome` sinon (bien sur, `xxx` est ici le palindrome en question). Pensez à vous débarrasser au préalable des majuscules et des espaces.

Testez ensuite si les expressions suivantes sont des palindromes :

- Radar
- Never odd or even
- Karine alla en Iran
- Un roc si biscornu

10.7.6 Mot composable

Un mot est composable à partir d'une séquence de lettres si la séquence contient toutes les lettres du mot. Comme au Scrabble, chaque lettre de la séquence ne peut être utilisée qu'une seule fois. Par exemple, `coucou` est composable à partir de `uocuoceokzefhu`.

Écrivez la fonction `composable()` qui prend en argument un mot (chaîne de caractères) et une séquence de lettre (chaîne de caractères) et qui affiche `Le mot xxx est composable à partir de yyy` si le mot (`xxx`) est composable à partir de la séquence de lettres (`yyy`) et `Le mot xxx n'est pas composable à partir de yyy` sinon.

Testez la fonction avec différents mots et séquences.

10.7.7 Alphabet et pangramme

Les codes ASCII des lettres minuscules de l'alphabet vont de 97 (lettre 'a') à 122 (lettre 'z'). La fonction `chr()` prend en argument un code ASCII sous forme d'un entier et renvoie le caractère correspondant. Ainsi `chr(97)` renvoie 'a', `chr(98)` renvoie 'b' et ainsi de suite.

Écrivez la fonction `get_alphabet()` qui utilise une boucle et la fonction `chr()` et qui renvoie une chaîne de caractères contenant toutes les lettres de l'alphabet.

Un pangramme⁴ est une phrase comportant au moins une fois chaque lettre de l'alphabet. Par exemple, "Portez ce vieux whisky au juge blond qui fume" est un pangramme.

3. http://en.wikipedia.org/wiki/Hamming_distance

4. <http://fr.wikipedia.org/wiki/Pangramme>

Écrivez la fonction `pangramme()` qui utilise la fonction `get_alphabet()` précédente, qui prend en argument une chaîne de caractère (`xxx`) et qui renvoie `xxx est un pangramme` si cette chaîne de caractères est un pangramme et `xxx n'est pas un pangramme` sinon. Pensez à vous débarrasser des majuscules le cas échéant.

Testez ensuite si les expressions suivantes sont des pangrammes :

- “Portez ce vieux whisky au juge blond qui fume”
- “Monsieur Jack vous dactylographiez bien mieux que votre ami Wolf”,
- “Buvez de ce whisky que le patron juge fameux”.

10.7.8 Affichage des carbones alpha d’une structure de protéine

Téléchargez le fichier `1bta.pdb` qui correspond à la structure tridimensionnelle de la protéine barstar⁵ sur le site de la PDB (lien direct vers le fichier⁶).

Écrivez la fonction `trouve_calpha()` qui prend en argument le nom d’un fichier PDB (sous forme de chaîne de caractères), qui sélectionne uniquement les lignes contenant des carbones alpha, qui stocke ces lignes et les renvoie sous forme de liste.

En utilisant la fonction `trouve_calpha()`, affichez à l’écran les carbones alpha des deux premiers résidus.

Mémo sur le format PDB

Dans un fichier au format PDB, les atomes qui compose une macromolécules sont présentés sous cette forme :

1	ATOM	681	N	GLY	A	43	-13.262	-2.363	-1.452	1.00	0.27		N
2	ATOM	682	CA	GLY	A	43	-14.258	-3.210	-0.729	1.00	0.31		C
3	ATOM	683	C	GLY	A	43	-13.627	-3.900	0.484	1.00	0.27		C
4	ATOM	684	O	GLY	A	43	-13.610	-5.112	0.574	1.00	0.32		O
5	ATOM	685	H	GLY	A	43	-12.322	-2.332	-1.164	1.00	0.37		H
6	ATOM	686	HA2	GLY	A	43	-15.075	-2.588	-0.396	1.00	0.39		H
7	ATOM	687	HA3	GLY	A	43	-14.638	-3.963	-1.404	1.00	0.37		H
8	ATOM	688	N	TRP	A	44	-13.125	-3.148	1.425	1.00	0.25		N
9	ATOM	689	CA	TRP	A	44	-12.519	-3.775	2.630	1.00	0.23		C
10	ATOM	690	C	TRP	A	44	-11.944	-2.689	3.540	1.00	0.23		C
11	ATOM	691	O	TRP	A	44	-12.415	-2.466	4.637	1.00	0.27		O

Pour sélectionner un carbone alpha, le champ délimité entre les caractères 1 et 6 doit correspondre à `ATOM` et celui délimité entre les caractères 13 et 16 doit correspondre à `CA`.

Les coordonnées tridimensionnelles (x, y et z) des atomes se trouvent dans les champs délimités par respectivement les caractères 31 à 38, 39 à 46 et 47 à 54.

N’oubliez pas que Python commence à compter à 0.

10.7.9 Calcul des distances entre les carbones alpha consécutifs d’une structure de protéine

En utilisant la fonction `trouve_calpha()` précédente, calculez la distance inter-atomique entre les carbones alpha des deux premiers résidus (avec deux chiffres après la virgule).

5. <http://www.rcsb.org/pdb/explore.do?structureId=1BTA>

6. <https://files.rcsb.org/download/1BTA.pdb>

Écrivez ensuite la fonction `calcule_distance()` qui prend en argument la liste renvoyée par la fonction `trouve_alpha()` précédente, qui calcule les distance inter-atomiques entre carbones alpha consécutifs et affiche ces distances sous la forme :

```
1 numero_calpha_1 numero_calpha_2 distance
```

La distance sera affichée avec deux chiffres après la virgule. Voici un exemple avec les premiers carbones alpha :

```
1 1 2 3.80
2 2 3 3.80
3 3 4 3.83
4 4 5 3.82
```

On rappelle que la distance d entre deux points A et B de coordonnées respectives (x_A, y_A, z_A) et (x_B, y_B, z_B) se calcule comme :

$$d = \sqrt{(x_B - x_A)^2 + (y_B - y_A)^2 + (z_B - z_A)^2}$$

Modifiez maintenant la fonction `calcule_distance()` pour qu'elle affiche à la fin la moyenne des distances.

La distance inter carbones alpha dans les protéines est très stable et de l'ordre de 3,8 angströms. Observez avec attention les valeurs que vous avez calculées pour la protéine 1BTA. Expliquez la valeur surprenante que vous devriez avoir obtenue.

Chapitre 11

Plus sur les listes

11.1 Propriétés des listes

Comme pour les chaînes de caractères, les listes possèdent de nombreuses **méthodes** (on rappelle, une *méthode* est une fonction qui agit sur l'objet à laquelle elle est attachée par un `.`) qui leur sont propres et qui peuvent se révéler très pratiques. Observez les exemples suivants :

- `append()` que nous avons déjà vue au chapitre 4 et qui permet d'ajouter un élément à la fin d'une liste existante.

```
1 >>> a = [1,2,3]
2 >>> a.append(5)
3 >>> a
4 [1, 2, 3, 5]
```

qui est équivalent à

```
1 >>> a = [1,2,3]
2 >>> a = 1 + [5]
3 >>> a
4 [1, 2, 3, 5]
```

Conseil : préférez la version avec `append()` qui est plus compacte et facile à lire.

- `insert()` pour insérer un objet dans une liste avec un indice déterminé.

```
1 >>> a.insert(2,-15)
2 >>> a
3 [1, 2, -15, 3, 5]
```

- `del` pour supprimer un élément d'une liste à une indice déterminé.

```
1 >>> del a[1]
2 >>> a
3 [1, -15, 3, 5]
```

Remarque : contrairement aux autres méthodes associées aux listes, `del` est une fonction générale de Python (utilisable pour d'autres objets que les listes), et celle-ci ne prend pas de parenthèse.

- `remove()` pour supprimer un élément d'une liste à partir de sa valeur.

```
1 >>> a.remove(5)
2 >>> a
3 [1, -15, 3]
```

— `sort()` pour trier une liste.

```
1 >>> a.sort()
2 >>> a
3 [-15, 1, 3]
```

— `reverse()` pour inverser une liste.

```
1 >>> a.reverse()
2 >>> a
3 [3, 1, -15]
```

— `count()` pour compter le nombre d'éléments (passé en argument) dans une liste.

```
1 >>> a=[1, 2, 4, 3, 1, 1]
2 >>> a.count(1)
3 3
4 >>> a.count(4)
5 1
6 >>> a.count(23)
7 0
```

Remarque 1 : attention, dans de nombreux exemples de méthodes ci-dessus, une liste remaniée (par exemple `a.sort()`) n'est pas renvoyée (la liste est modifiée de manière interne, mais l'appel de la méthode ne renvoie rien, c'est à dire, pas d'objet récupérable dans une variable) ! Il s'agit d'un exemple d'utilisation de méthode (donc de fonction) qui fait une action mais qui ne renvoie rien. Pensez-y dans vos utilisations futures des listes.

Remarque 2 : attention, certaines fonctions ci-dessus décalent les indices d'une liste (par exemple `insert()`, `del` etc).

La méthode `append()`, que nous avons déjà vue, est particulièrement pratique car elle permet de construire une liste au fur et à mesure des itérations d'une boucle. Pour cela, nous vous rappelons qu'il est commode de définir préalablement une liste vide de la forme `ma_liste = []`. Voici un exemple où une chaîne de caractères est convertie en liste :

```
1 >>> seq = 'CAAAGGTAACGC'
2 >>> seq_list = []
3 >>> seq_list
4 []
5 >>> for base in seq:
6 ...     seq_list.append(base)
7 ...
8 >>> seq_list
9 ['C', 'A', 'A', 'A', 'G', 'G', 'T', 'A', 'A', 'C', 'G', 'C']
```

Remarquez que vous pouvez directement utiliser la fonction `list()` qui prend n'importe quel objet séquentiel (liste, chaîne de caractères, *tuples* [que nous verrons au chapitre 13], etc.) et qui renvoie une liste :

```
1 >>> seq = 'CAAAGGTAACGC'
2 >>> list(seq)
3 ['C', 'A', 'A', 'A', 'G', 'G', 'T', 'A', 'A', 'C', 'G', 'C']
```

Cette méthode est certes plus simple, mais il arrive parfois que l'on doive utiliser les boucles tout de même, comme lorsqu'on lit un fichier. On rappelle l'action `list(seq)` convertit un objet de type chaîne

de caractères en un objet de type liste (il s'agit donc de **casting**). De même que `list(range(10))` convertissait un objet de type `range` en un objet de type liste.

Enfin, si vous voulez avoir une liste exhaustive des méthodes disponibles, nous rappelons que vous pouvez utiliser la fonction `dir(l)` (`l` étant un objet de type `list`).

11.2 Test d'appartenance

L'opérateur `in` permet de tester si un élément fait partie d'une liste.

```

1 liste = [1, 3, 5, 7, 9]
2 >>> 3 in liste
3 True
4 >>> 4 in liste
5 False
6 >>> 3 not in liste
7 False
8 >>> 4 not in liste
9 True

```

La variation avec `not` permet, *a contrario*, de vérifier qu'un élément n'est pas dans une liste.

11.3 Copie de listes

Il est très important de savoir que l'affectation d'une liste (à partir d'une liste préexistante) crée en réalité une **référence** et non une **copie** :

```

1 >>> x = [1,2,3]
2 >>> y = x
3 >>> y
4 [1, 2, 3]
5 >>> x[1] = -15
6 >>> x
7 [1, -15, 3]
8 >>> y
9 [1, -15, 3]

```

Vous voyez que la modification de `x` modifie `y` aussi ! Pour comprendre ce qu'il se passe nous allons de nouveau utiliser le site [pythontutor](http://www.pythontutor.com/)¹ sur cet exemple :

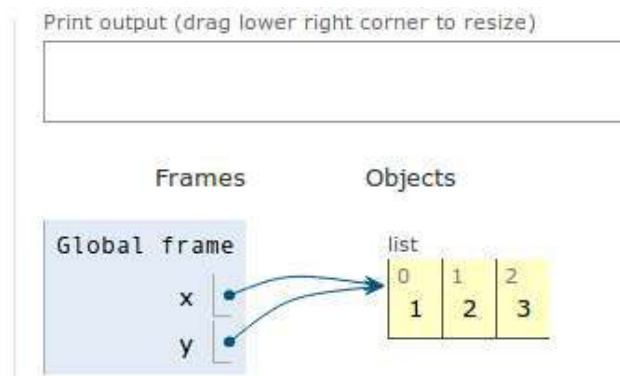
Python 3.6

```

1 x = [1,2,3]
2 y = x
3 print(y)
4 x[1] = -15
5 print(y)

```

[t code](#) | [Live programming](#)



1. <http://www.pythontutor.com/>

Techniquement, Python utilise des pointeurs (comme dans le langage C) vers les mêmes objets. *Pythontutor* nous l'illustre en mettant des flèches, et on voit bien que `x` et `y` pointent vers la même liste. Donc si on modifie `x`, `y` est modifiée pareillement. Rappelez-vous de ceci dans vos futurs programmes car cela pourrait avoir des effets désastreux !

Pour éviter le problème, il va falloir créer une copie explicite de la liste initiale. Regardez cet exemple :

```

1 >>> x = [1,2,3]
2 >>> y = x[:]
3 >>> x[1] = -15
4 >>> y
5 [1, 2, 3]
```

Dans l'exemple précédent, `x[:]` a créé une copie *à la volée* de la liste `x`. Vous pouvez utiliser aussi la fonction `list()` qui renvoie explicitement une liste :

```

1 >>> x = [1,2,3]
2 >>> y = list(x)
3 >>> x[1] = -15
4 >>> y
5 [1, 2, 3]
```

Si on regarde à nouveau dans *pythontutor*, on voit bien que l'utilisation des tranches `[:]` ou de la fonction `list()` crée des copies explicites, c'est à dire que chaque "flèche" pointe vers des listes différentes, indépendantes les unes des autres.

Python 3.6

```

1 x = [1, 2, 3]
2 y = x[:]
3 z = list(x)
4 x[1] = -15
5 print(y)
6 print(z)
```

[Edit code](#) | [Live programming](#)

ited

breakpoint; use the Back and Forward buttons to jump there.

< Back Step 5 of 6 Forward > Last >>

Print output (drag lower right corner to resize)

Frames	Objects																		
<div style="border: 1px solid gray; padding: 5px;"> <p>Global frame</p> <ul style="list-style-type: none"> x y z </div>	<div style="margin-bottom: 10px;"> <p>list</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td>0</td><td>1</td><td>2</td></tr> <tr><td>1</td><td>-15</td><td>3</td></tr> </table> </div> <div style="margin-bottom: 10px;"> <p>list</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td>0</td><td>1</td><td>2</td></tr> <tr><td>1</td><td>2</td><td>3</td></tr> </table> </div> <div> <p>list</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td>0</td><td>1</td><td>2</td></tr> <tr><td>1</td><td>2</td><td>3</td></tr> </table> </div>	0	1	2	1	-15	3	0	1	2	1	2	3	0	1	2	1	2	3
0	1	2																	
1	-15	3																	
0	1	2																	
1	2	3																	
0	1	2																	
1	2	3																	

Attention, les deux techniques précédentes ne fonctionnent que pour les listes à une dimension, autrement dit les listes qui ne contiennent pas elles-mêmes d'autres listes. Voyez par exemple :

```

1 >>> x = [[1,2],[3,4]]
2 >>> x
3 [[1, 2], [3, 4]]
4 >>> y = x[:]
5 >>> x[1][1] = 55
6 >>> x
7 [[1, 2], [3, 55]]
8 >>> y
```

```
9 [[1, 2], [3, 55]]
```

et

```
1 >>> y = list(x)
2 >>> x[1][1] = 77
3 >>> x
4 [[1, 2], [3, 77]]
5 >>> y
6 [[1, 2], [3, 77]]
```

La méthode de copie qui **fonctionne à tous les coups** consiste à appeler la fonction `deepcopy()` du module `copy`.

```
1 >>> import copy
2 >>> x = [[1,2],[3,4]]
3 >>> x
4 [[1, 2], [3, 4]]
5 >>> y = copy.deepcopy(x)
6 >>> x[1][1] = 99
7 >>> x
8 [[1, 2], [3, 99]]
9 >>> y
10 [[1, 2], [3, 4]]
```

11.4 Exercices

Conseil : pour ces exercices, écrivez des scripts dans des fichiers, puis exécutez-les dans un *shell*.

11.4.1 Tri de liste

Soit la liste de nombres `[8, 3, 12.5, 45, 25.5, 52, 1]`. Triez les nombres de cette liste par ordre croissant, sans utiliser la fonction `sort()` (les fonctions/méthodes `min()`, `append()` et `remove()` vous seront utiles).

11.4.2 Séquence nucléique aléatoire

Générez aléatoirement une séquence nucléique de 20 bases en utilisant une liste et la méthode `append()` (chaque base est équiprobable). Votre programme contiendra une fonction `gen_seq_aléa()` à laquelle on passe un entier positif `n` représentant le nombre de nucléotides. La fonction renverra la séquence aléatoire sous forme de liste.

11.4.3 Séquence nucléique complémentaire

Reprenez l'exercice 9.5.5. Transformez la séquence nucléique `TCTGTTAACCATCCACTTCG` en sa séquence complémentaire inverse en utilisant toutes un maximum de méthodes associées aux listes. N'oubliez pas que la séquence complémentaire doit être inversée, pensez aux méthodes des listes !

11.4.4 Doublons

Soit la liste de nombres [5, 1, 1, 2, 5, 6, 3, 4, 4, 4, 2]. Enlevez les doublons de cette liste, triez-là et affichez-là.

11.4.5 Séquence nucléique aléatoire 2

Générez aléatoirement une séquence nucléique de 50 bases contenant 10 % de A, 50 % de G, 30 % de T et 10 % de C. Votre programme contiendra une fonction `gen_seq_alea2()` à qui on passe 4 *floats* représentant les pourcentages de chaque nucléotide. La fonction renverra la séquence sous forme de liste.

11.4.6 Triangle de Pascal (Exercice +++)

Voici le début du triangle de Pascal :

```
1 1
2 11
3 121
4 1331
5 14641
6 ...
```

Comprenez comment une ligne est construite à partir de la précédente. A partir de l'ordre 1 (ligne 2, 11), générez l'ordre suivant (121). Vous pouvez utiliser une liste préalablement générée avec `range()`. Généralisez à l'aide d'une boucle. Ecrivez dans un fichier `pascal.out` les lignes du triangle de Pascal de l'ordre 1 jusqu'à l'ordre 10.

Chapitre 12

Plus sur les fonctions

Avant d'attaquer ce chapitre, nous vous conseillons de relire le chapitre 9 et de bien en assimiler toutes les notions (et aussi en faire les exercices). Nous avons vu dans ce chapitre 9 le concept puissant et incontournable que représente les **fonctions**. Nous vous avons également introduit la notion de variables **locales** et **globales**. Dans ce chapitre, nous allons aller un peu plus loin sur la visibilité de ces variables dans et hors des fonction, et aussi voir ce qu'il se passe lorsque ces variables sont des listes. Attention, la plupart des lignes de code ci-dessous sont données à titre d'exemple pour bien comprendre ce qu'il se passe, mais nombre d'entre elles sont des aberrations de programmation. Ainsi nous ferons un récapitulatif des bonnes pratiques à la fin du chapitre. Enfin, nous vous conseillons de tester tous les exemples ci-dessous avec l'excellent site [pythontutor](http://www.pythontutor.com/)¹ afin de suivre l'état des variables lors de l'exécution des exemples.

12.1 Appel d'une fonction dans une fonction

Dans le chapitre 9 nous avons vu des fonctions qui étaient appelées depuis le programme principal. Il est en fait possible d'appeler une fonction depuis une autre fonction (et plus généralement on peut appeler une fonction de n'importe où si tant est qu'elle est visible par Python, *i.e.* chargée dans la mémoire). Regardez cet exemple :

```
1 # définition des fonctions
2 def polynome(x):
3     return (x**2 - 2*x + 1)
4
5 def calc_vals(debut, fin):
6     liste_vals = []
7     for x in range(debut, fin + 1):
8         liste_vals.append(polynome(x))
9     return liste_vals
10
11 # programme principal
12 print(calc_vals(-5,5))
```

Dans cet exemple, nous appelons depuis le programme principal la fonction `calc_vals()`, puis à l'intérieur de celle-ci nous appelons l'autre fonction `polynome()`. Regardons ce que *pythontutor* nous montre lorsque la fonction `polynome()` est en exécution (lors de son second appel) :

1. <http://www.pythontutor.com/>

Python 3.6

```

1 def polynome(x):
2     return (x**2 - 2*x + 1)
3
4 def calc_vals(debut, fin):
5     liste_vals = []
6     for x in range(debut, fin + 1):
7         liste_vals.append(polynome(x))
8     return liste_vals
9
10 # programme principal
11 print(calc_vals(-5,5))

```

[Edit this code](#)

at has just executed
 re to execute
 of code to set a breakpoint; use the Back and Forward buttons to jump there.

<< First < Back Step 15 of 63 Forward > Last >>

Print output (drag lower right corner to resize)

Frames	Objects
Global frame	function polynome(x)
polynome	function calc_vals(debut, fin)
calc_vals	list
debut	-5
fin	5
liste_vals	[0, 36]
x	-4
polynome	x: -4
Return value	25

Nous voyons que l'espace mémoire alloué à `polynome()` est bleuté, nous indiquant que cette fonction est en cours d'exécution. La fonction appelante `calc_vals()` est toujours là (sur un fond blanc) car son exécution n'est pas terminée. Elle est en quelque sorte *figée* dans le même état qu'avant l'appel de `polynome()`, et on pourra ainsi noter que ses variables *locales* sont toujours là (c'est à dire `debut`, `fin`, `liste_vals` et `x`). De manière générale, les variables *locales* d'une fonction ne seront détruites que lorsque l'exécution de celle-ci sera terminée. Dans notre exemple, les variables *locales* de `calc_vals()` ne seront détruites que lorsque la boucle sera terminée et que la liste `liste_vals` sera retournée au programme principal. En passant, la fonction `calc_vals()` appelle la fonction `polynome()` à chaque itération de la boucle.

Ainsi le programmeur est libre de faire tous les appels qu'il souhaite. Une fonction peut appeler une autre fonction, cette dernière peut appeler une autre fonction et ainsi de suite (et autant de fois qu'on le veut). Une fonction peut même s'appeler elle-même, cela s'appelle une fonction *récursive* (voir la section suivante). Attention toutefois à retrouver vos petits si vous vous perdez dans les appels successifs !

12.2 Fonctions récursives

Conseil : Pour les débutants, vous pouvez passer cette section.

Comme dit ci-dessus, une fonction récursive est une fonction qui s'appelle elle-même. Les fonctions récursives permettent d'obtenir une efficacité redoutable dans la résolution de certains algorithmes comme le tri rapide² (en anglais *quicksort*). Oublions la recherche d'efficacité pour l'instant, et concentrons nous sur un exemple de base, la fonction factorielle, qui illustre à merveille comment sont gérés des fonctions récursives en Python. Regardez ce code et tenter de comprendre ce qu'il va faire :

```

1 def calc_factorielle(nb):
2     if nb == 1:
3         return 1
4     else:
5         return nb * calc_factorielle(nb - 1)

```

2. https://fr.wikipedia.org/wiki/Tri_rapide

```

6
7 # prog principal
8 print(calc_factorielle(4))

```

Pas très facile, n'est-ce pas ? A nouveau, demandons l'aide de pythontutor pour voir ce qu'il se passe (nous vous conseillons bien sûr de tester vous même cet exemple sur pythontutor³) :

Python 3.6

```

1 def calc_factorielle(nb):
2     if nb == 1:
3         return 1
4     else:
5         return nb * calc_factorielle(nb - 1)
6
7 # prog principal
8 print(calc_factorielle(4))

```

[Edit this code](#)

that has just executed
t line to execute
line of code to set a breakpoint; use the Back and Forward buttons to jump there.

Step 15 of 18

Help improve this tool by clicking whenever you learn something:

Print output (drag lower right corner to resize)

Frames	Objects
Global frame	function calc_factorielle(nb)
calc_factorielle	nb 4
calc_factorielle	nb 3
calc_factorielle	nb 2
calc_factorielle	nb 1 Return value 1

Vous voyez qu'au début on appelle la fonction en passant la variable `nb` qui vaut 4. Puis au sein de la fonction, celle-ci se rappelle elle même, mais cette fois-ci en passant la valeur 3. Au prochain appel, ce sera avec la valeur 2, puis finalement 1. Dans ce dernier cas, le `if nb == 1:` est vrai, et la ligne `return 1` va être exécutée. A ce moment précis de l'exécution, vous voyez que les appels successifs forment une sorte de *pile* (voir image ci dessus [AJOUTER ICI UN LIEN]). La valeur 1 va être ainsi renvoyée au niveau de l'appel précédent, puis le résultat $2 * 1 = 2$ (où 2 correspond à `nb` et 1 provient de `calc_factorielle(nb - 1) = 1`) va être renvoyé à l'appel précédent, puis $3 * 2 = 6$ (où 3 correspond à `nb` et 2 provient de `calc_factorielle(nb - 1) = 2`) va être renvoyé à l'appel précédent, pour finir par $4 * 6 = 24$ (où 4 correspond à `nb` et 6 provient de `calc_factorielle(nb - 1) = 6`), soit la valeur de `!`4. Les appels successifs vont donc se "dépiler" et nous reviendrons dans le programme principal.

Même si les fonctions récursives peuvent être ardues à comprendre, notre propos est ici de vous illustrer qu'une fonction qui en appelle une autre (ici il s'agit d'elle même) reste "figée" dans le même état, jusqu'à temps que la fonction appelée lui renvoie une valeur.

12.3 Portée des variables

Il est très important lorsque l'on manipule des fonctions de connaître la portée des variables. On a vu que les variables créées au sein d'une fonction ne sont pas visibles à l'extérieur de celle-ci car elles étaient **locales** à la fonction. Observez le code suivant :

```

1 >>> def mafonction():

```

3. <http://pythontutor.com>

```

2 ...     x = 2
3 ...     print('x vaut {} dans la fonction'.format(x))
4 ...
5 >>> mafonction()
6 x vaut 2 dans la fonction
7 >>> print(x)
8 Traceback (most recent call last):
9   File "<stdin>", line 1, in <module>
10 NameError: name 'x' is not defined

```

Lorsque Python exécute le code de la fonction, il connaît le contenu de la variable `x`. Par contre, de retour dans le module principal (dans notre cas, il s'agit de l'interpréteur Python), il ne la connaît plus d'où le message d'erreur.

De même, une variable passée en argument est considérée comme **locale** lorsqu'on arrive dans la fonction :

```

1 >>> def mafonction(x):
2 ...     print('x vaut {} dans la fonction'.format(x))
3 ...
4 >>> mafonction(2)
5 x vaut 2 dans la fonction
6 >>> print(x)
7 Traceback (most recent call last):
8   File "<stdin>", line 1, in ?
9 NameError: name 'x' is not defined

```

Deuxièmement, lorsqu'une variable déclarée à la racine du module (c'est comme cela que l'on appelle un programme Python), elle est visible dans tout le module. On a vu qu'on parlait de variable **globale**

```

1 >>> def mafonction():
2 ...     print(x)
3 ...
4 >>> x = 3
5 >>> mafonction()
6 3
7 >>> print(x)
8 3

```

Dans ce cas, la variable `x` est visible dans le module principal et dans toutes les fonctions du module. Toutefois, Python ne permet pas la modification d'une variable globale dans une fonction :

```

1 >>> def mafonction():
2 ...     x = x + 1
3 ...
4 >>> x = 1
5 >>> mafonction()
6 Traceback (most recent call last):
7   File "<stdin>", line 1, in <module>
8     File "<stdin>", line 2, in fct
9 UnboundLocalError: local variable 'x' referenced before assignment

```

L'erreur renvoyée montre que Python pense que `x` est une variable locale qui n'a pas été encore assignée. Si on veut vraiment modifier une variable globale dans une fonction, il faut utiliser le mot-clé `global` :


```

1 >>> def mafonction():
2 ...     global x
3 ...     x = x + 1
4 ...
5 >>> x = 1
6 >>> mafonction()
7 >>> x
8 2

```

Dans ce dernier cas, le mot-clé `global` a forcé la variable `x` à être globale plutôt que locale au sein de la fonction.

12.4 Portée des listes

Préambule : Attention, les exemples de cette section représentent des absurdités de programmation. Ils sont donnés à titre indicatif pour comprendre ce qu'il se passe, mais il ne faut surtout pas s'en inspirer !

Soyez extrêmement attentifs avec les types modifiables (tels que les listes) car vous pouvez les changer au sein d'une fonction :

```

1 >>> def mafonction():
2 ...     liste[1] = -127
3 ...
4 >>> liste = [1,2,3]
5 >>> mafonction()
6 >>> liste
7 [1, -127, 3]

```

De même que si vous passez une liste en argument, elle est tout autant modifiable au sein de la fonction :

```

1 >>> def mafonction(x):
2 ...     x[1] = -15
3 ...
4 >>> y = [1,2,3]
5 >>> mafonction(y)
6 >>> y
7 [1, -15, 3]

```

Si vous voulez éviter ce problème, utilisez des tuples, Python renverra une erreur puisque ces derniers sont non modifiables ! Une autre solution pour éviter la modification d'une liste lorsqu'elle est passée en tant qu'argument, est de la passer explicitement (comme nous l'avons fait pour la copie de liste) afin qu'elle reste intacte dans le programme principal.

```

1 >>> def mafonction(x):
2 ...     x[1] = -15
3 ...
4 >>> y = [1,2,3]
5 >>> mafonction(y[:])
6 >>> y
7 [1, 2, 3]
8 >>> mafonction(list(y))

```

```

9 >>> y
10 [1, 2, 3]

```

Dans ces deux derniers exemples, une copie de `y` est créée à la volée lorsqu'on appelle la fonction, ainsi la liste `y` du module principal reste intacte.

12.5 Règle LGI

Lorsque Python rencontre une variable, il va traiter la résolution de son nom avec des priorités particulières : d'abord il va regarder si la variable est **locale**, puis si elle n'existe pas localement, il vérifiera si elle est **globale** et enfin si elle n'est pas globale, il testera si elle est **interne** (par exemple la fonction `len()` est considérée comme une fonction interne à Python, *i.e.* elle existe à chaque fois que vous lancez Python). On appelle cette règle la règle **LGI** pour locale, globale, interne. En voici un exemple :

```

1 >>> def mafonction():
2 ...     x = 4
3 ...     print('Dans la fonction x vaut {}'.format(x))
4 ...
5 >>> x = -15
6 >>> mafonction()
7 Dans la fonction x vaut 4
8 >>> print('Dans le module principal x vaut {}'.format(x))
9 Dans le module principal x vaut -15

```

Vous voyez que dans la fonction, `x` a pris la valeur qui lui était définie localement en priorité sur sa valeur définie dans le module principal.

Conseil : même si Python peut reconnaître une variable ayant le même nom que ses fonctions ou variables internes, évitez de les utiliser car ceci rendra votre code confus !

12.6 Recommandations

Dans ce chapitre nous avons *joué* avec les fonctions (et les listes) afin de vous montrer comment Python réagit. Toutefois, notez bien que **l'utilisation de variables globales est à banir définitivement de votre pratique de la programmation**. Parfois on se dit, je n'ai pas le temps, je préfère créer une variable globale visible partout dans le programme (donc dans toutes les fonctions), car "*ça va plus vite, c'est plus simple*". C'est un très mauvais calcul, ne serait ce que parce-que vos fonctions ne seront pas réutilisables dans un autre contexte si elles utilisent des variables globales ! Ensuite, arriverez-vous à vous relire dans 6 mois ? Et si vous donnez votre code à lire à quelqu'un ? Il existe de nombreuses autres raisons⁴ que nous ne développerons pas ici, mais libre à vous de consulter de la documentation externe.

Heureusement, Python est orienté objet et cela permet "*d'encapsuler*" des variables dans des objets et de s'affranchir définitivement des variables globales (nous verrons cela dans le chapitre 18). En attendant, et si vous souhaitez ne pas aller plus loin sur les notions d'objet (on peut tout à fait *pythonner* sans cela), retenez la chose suivante avec les fonctions et les variables globales :

Plutôt que d'utiliser des variables globales, passez vos variables explicitement aux fonctions en tant qu'argument(s).

4. <http://wiki.c2.com/?GlobalVariablesAreBad>

Nous espérons que vous maîtrisez maintenant les fonctions sous tous leurs angles. Comme dit en introduction du chapitre 9, elles sont incontournables et tout programmeur se doit de les maîtriser. Voici les derniers conseils que nous pouvons vous donner :

- Lorsque vous attaquez un nouveau projet de programmation complexe, posez-vous la question : *“Comment pourriez-vous décomposer en blocs chaque tâche à effectuer, chaque bloc pouvant être une fonction”*. Si une fonction s’avère trop complexe, vous pouvez également la décomposer en d’autres fonctions.
- Au risque de nous répéter, forcez vous à utiliser les fonctions en permanence, pratiquez, pratiquez... et pratiquez encore !

12.7 Exercices

Conseil : pour ces exercices, écrivez des scripts dans des fichiers, puis exécutez-les dans un *shell*.

12.7.1 Prédire la sortie

- Prédisez le comportement de ce code sans le recopier dans un script ni dans l’interpréteur Python :

```

1 def hello(prenom):
2     print("Bonjour {}".format(prenom))
3
4 hello("Patrick")
5 print(x)
```

- Prédisez le comportement de ce code sans le recopier dans un script ni dans l’interpréteur Python :

```

1 x = 10
2
3 def hello(prenom):
4     print("Bonjour {}".format(prenom))
5
6 hello("Patrick")
7 print(x)
```

- Prédisez le comportement de ce code sans le recopier dans un script ni dans l’interpréteur Python :

```

1 x = 10
2
3 def hello(prenom):
4     print("Bonjour {}".format(prenom))
5     print(x)
6
7 hello("Patrick")
8 print(x)
```

- Prédisez le comportement de ce code sans le recopier dans un script ni dans l’interpréteur Python :

```

1 x = 10
2
3 def hello(prenom):
```

```
4 | x = 42
5 | print("Bonjour {}".format(prenom))
6 | print(x)
7 |
8 | hello("Patrick")
9 | print(x)
```

Chapitre 13

Dictionnaires et tuples

13.1 Dictionnaires

Les **dictionnaires** se révèlent très pratiques lorsque vous devez manipuler des structures complexes à décrire et que les listes présentent leurs limites. Les dictionnaires sont des collections non ordonnées d'objets, c'est à dire qu'il n'y a pas de notion d'ordre (*i.e.* pas d'indice). On accède aux **valeurs** d'un dictionnaire par des **clés**. Ceci semble un peu confus? Regardez l'exemple suivant :

```
1 >>> ani1 = {}
2 >>> ani1['nom'] = 'girafe'
3 >>> ani1['taille'] = 5.0
4 >>> ani1['poids'] = 1100
5 >>> ani1
6 {'nom': 'girafe', 'poids': 1100, 'taille': 5.0}
```

En premier, on définit un dictionnaire vide avec les symboles {} (tout comme on peut le faire pour les listes avec []). Ensuite, on remplit le dictionnaire avec différentes clés ('nom', 'taille', 'poids') auxquelles on affecte des valeurs ('girafe', 5.0, 1100). Vous pouvez mettre autant de clés que vous voulez dans un dictionnaire (tout comme vous pouvez ajouter autant d'éléments que vous voulez dans une liste).

On peut aussi initialiser toutes les clés et les valeurs d'un dictionnaire en une seule opération :

```
1 >>> ani2 = {'nom':'singe', 'poids':70, 'taille':1.75}
```

Mais rien ne nous empêche d'ajouter une clé et une valeur supplémentaire :

```
1 >>> ani2[age] = 15
```

Pour récupérer la valeur associée à une clé donnée, il suffit d'utiliser la syntaxe `dictionnaire['cle']` :

```
1 >>> ani1['taille']
2 5.0
```

13.1.1 Itération sur les clefs pour obtenir les valeurs

Il est possible d'obtenir toutes les valeurs d'un dictionnaire à partir de ses clés :

```
1 >>> ani2 = {'nom':'singe', 'poids':70, 'taille':1.75}
2 >>> for clef in ani2:
```

```
3 ...     print(clef, ani2[clef])
4 ...
5 poids 70
6 nom singe
7 taille 1.75
```

13.1.2 Méthodes `keys()` et `values()`

Les méthodes `.keys()` et `.values()` renvoient, comme vous vous en doutez, les clés et les valeurs d'un dictionnaire :

```
1 >>> ani2.keys()
2 dict_keys(['poids', 'nom', 'taille'])
3 >>> ani2.values()
4 dict_values([70, 'singe', 1.75])
```

Les mentions `dict_keys` et `dict_values` indiquent que nous avons à faire à des objets un peu particulier. Si besoin, nous pouvons les transformer en liste avec la fonction `list()` (par exemple : `list(ani2.values())`).

13.1.3 Existence d'une clef

Pour vérifier si une clé existe dans un dictionnaire :

```
1 >>> if "poids" in ani2:
2 ...     print("La clef 'poids' existe pour ani2")
3 ...
4 La clef 'poids' existe pour ani2
```

13.1.4 Liste de dictionnaires

En créant une liste de dictionnaires qui possèdent les mêmes clés, on obtient une structure qui ressemble à une base de données :

```
1 >>> animaux = [ani1, ani2]
2 >>> animaux
3 [{'nom': 'girafe', 'poids': 1100, 'taille': 5.0}, {'nom': 'singe',
4 'poids': 70, 'taille': 1.75}]
5 >>>
6 >>> for ani in animaux:
7 ...     print(ani['nom'])
8 ...
9 girafe
10 singe
```

Vous voyez que les dictionnaires permettent de gérer des structures complexes de manière plus explicite que les listes.

13.2 Tuples

Les **tuples** correspondent aux listes à la différence qu'ils sont **non modifiables**. On a vu à la section précédente que les listes pouvaient être modifiées par des références ; les tuples vous permettent de vous affranchir de ce problème puisqu'ils sont non modifiables. Pratiquement, ils utilisent les parenthèses au lieu des crochets :

```

1 >>> x = (1,2,3)
2 >>> x
3 (1, 2, 3)
4 >>> x[2]
5 3
6 >>> x[0:2]
7 (1, 2)
8 >>> x[2] = 15
9 Traceback (innermost last):
10 File "<stdin>", line 1, in ?
11 TypeError: object doesn't support item assignment

```

L'affectation et l'indigage fonctionne comme avec les listes, mais si l'on essaie de modifier un des éléments du tuple, Python renvoie un message d'erreur. Si vous voulez ajouter un élément (ou le modifier), vous devez créer un autre tuple :

```

1 >>> x = (1,2,3)
2 >>> x + (2,)
3 (1, 2, 3, 2)

```

Remarque : pour utiliser un tuple d'un seul élément, vous devez utiliser une syntaxe avec une virgule (**element,**), ceci pour éviter une ambiguïté avec une simple expression.

Autre particularité des tuples, il est possible d'en créer de nouveaux sans les parenthèses, dès lors que ceci ne pose pas d'ambiguïté avec une autre expression :

```

1 >>> x = (1,2,3)
2 >>> x
3 (1, 2, 3)
4 >>> x = 1,2,3
5 >>> x
6 (1, 2, 3)

```

Toutefois, nous vous conseillons d'utiliser systématiquement les parenthèses afin d'éviter les confusions.

Enfin, on peut utiliser la fonction `tuple(sequence)` qui fonctionne exactement comme la fonction `list()`, c'est à dire qu'elle prend en argument un objet séquentiel et renvoie le tuple correspondant (opération de *casting*) :

```

1 >>> tuple([1,2,3])
2 (1, 2, 3)
3 >>> tuple("ATGCCGCGAT")
4 ('A', 'T', 'G', 'C', 'C', 'G', 'C', 'G', 'A', 'T')

```

Remarque : les listes, dictionnaires, tuples sont des objets qui peuvent contenir des collections d'autres objets. On peut donc construire des listes qui contiennent des dictionnaires, des tuples ou d'autres listes, mais aussi des dictionnaires contenant des tuples, des listes, etc.

Pratiquement, nous avons déjà croisé les *tuples* avec la fonction `enumerate()` (cf. chapitre 5) et aussi avec les fonctions lorsqu'on voulait renvoyer plusieurs valeurs (cf. chapitre 9, par exemple dans `return x,y`, `x,y` est un tuple).

13.3 Exercices

Conseil : pour ces exercices, écrivez des scripts dans des fichiers, puis exécutez-les dans un *shell*.

13.3.1 Composition en acides aminés

En utilisant un dictionnaire, déterminez le nombre d'occurrences de chaque acide aminé dans la séquence AGWPSGGASAGLAILWGASAIMPGALW. Le dictionnaire ne doit contenir que les acides aminés présents dans la séquence.

13.3.2 Mots de 2 lettres

Soit la séquence nucléotidique suivante :

```
1 ACCTAGCCATGTAGAATCGCCTAGGCTTTAGCTAGCTCTAGCTAGCTG
```

En utilisant un dictionnaire, faites un programme qui répertorie tous les mots de 2 lettres qui existent dans la séquence (AA, AC, AG, AT, etc.) ainsi que leur nombre d'occurrences puis qui les affiche à l'écran.

13.3.3 Mots de 3 et 4 lettres

Faites de même avec des mots de 3 et 4 lettres.

13.3.4 Mots de 2 lettres de *Saccharomyces cerevisiae*

En vous basant sur les scripts précédents, extrayez les mots de 2 lettres et leur occurrence sur le génome du chromosome I de la levure du boulanger *Saccharomyces cerevisiae* (fichier NC_001133.fna : lien GenBank¹, lien local²). Attention, le génome complet est fourni au format fasta.

13.3.5 Mots de n lettres et fichiers genbank

Créez un script `extract-words.py` qui prend en arguments un fichier genbank suivi d'un entier compris entre 1 et 4. Ce script doit extraire du fichier genbank tous les mots (ainsi que leur nombre d'occurrences) du nombre de lettres passées en option.

13.3.6 Mots de n lettres du génome d'*E. Coli*

Appliquez ce script sur le génome d'*Escherichia coli* : (fichier NC_000913.fna : lien GenBank³, lien local⁴). Attention, le génome complet est fourni au format fasta. Cette méthode vous paraît-elle efficace sur un génome assez gros comme celui d'*E. Coli*? Comment pourrait-on en améliorer la rapidité?

1. ftp://ftp.ncbi.nlm.nih.gov/genomes/Fungi/Saccharomyces_cerevisiae_uid128/NC_001133.fna

2. data-files/NC_001133.fna

3. ftp://ftp.ncbi.nlm.nih.gov/genomes/Bacteria/Escherichia_coli_K_12_substr__MG1655_uid57779/NC_000913.fna

4. data-files/NC_000913.fna

13.3.7 Dictionnaire et carbone alpha

À partir du fichier PDB 1BTA⁵, construisez un dictionnaire qui contient 4 clés se référant au premier carbone alpha : le numéro du résidu, puis les coordonnées atomiques x , y et z .

13.3.8 Dictionnaire et PDB

Sur le même modèle que ci-dessus, créez une liste de dictionnaires pour chacun des carbones alpha de la protéine.

13.3.9 Barycentre d'une protéine

À l'aide de cette liste, calculez les coordonnées x , y et z du barycentre de ces carbones alpha.

5. <http://www.rcsb.org/pdb/files/1BTA.pdb>

Chapitre 14

Création de modules

14.1 Pourquoi créer ses propres modules ?

Dans le chapitre 8, nous avons vu comment utiliser les modules existant dans python (comme `random`, `math`, etc). Nous avons vu par ailleurs les fonctions dans les chapitres 9 et 12 qui permettent de ré-utiliser une fraction de code plusieurs fois au sein d'un même programme sans avoir à dupliquer du code. On peut imaginer qu'une fonction bien écrite pourrait être judicieusement ré-utilisée dans un autre programme Python. C'est justement l'objet de la création de module, il s'agit d'y mettre un ensemble de fonctions que vous pourrez être amené à ré-utiliser souvent. En général les modules sont regroupés autour d'un thème précis, par exemple, on pourrait concevoir un module de gestion de séquences biologiques ou encore de gestion de fichiers PDB.

14.2 Comment créer son propre module

En Python, la création de modules est très simple. Il vous suffit d'écrire un ensemble de fonctions (et/ou de variables) dans un fichier, puis d'enregistrer ce dernier avec une extension `.py` (comme n'importe quel script Python). A titre d'exemple, nous allons créer un module simple que nous enregistrerons sous le nom `message.py`.

```
1 #!/usr/bin/env python3
2
3 """Module inutile qui affiche des messages :-)
4 """
5
6 def Bonjour(nom):
7     """Affiche bonjour !
8     """
9     return "bonjour " + nom
10
11 def Ciao(nom):
12     """Affiche ciao !
13     """
14     return "ciao " + nom
15
16 def Hello(nom):
17     """Affiche hello !
18     """
19     return "hello " + nom + " !"
```

```

20
21 date=16092008

```

Les commentaires entre triple guillemets en tête de module et dans chaque fonction sont facultatifs. On les appelle *docstrings* (voir section ci-dessous).

14.3 Utilisation de son propre module

Pour appeler une fonction ou une variable de ce module, il faut que le fichier `message.py` soit dans le répertoire courant (dans lequel on travaille) ou bien dans un répertoire indiqué par la variable d'environnement Unix `PYTHONPATH` (on rappelle en Unix il faut taper la commande suivante : `export PYTHONPATH=/chemin/vers/mes/super/modules/python`). Ensuite, il suffit d'importer le module et toutes ses fonctions (et variables) vous sont alors accessibles.

La première fois qu'un module est importé, Python crée un répertoire nommé `__pycache__` contenant un fichier avec une extension `.pyc` qui contient le bytecode¹ (code précompilé) du module.

L'importation du module se fait avec la commande `import message`. Notez que le fichier est bien enregistré avec une extension `.py` et pourtant on ne la précise pas lorsqu'on importe le module. Ensuite on peut utiliser les fonctions comme avec un module classique.

```

1 >>> import message
2 >>> message.Hello("Joe")
3 'hello Joe !'
4 >>> message.Ciao("Bill")
5 'ciao Bill'
6 >>> message.Bonjour("Monsieur")
7 'bonjour Monsieur'
8 >>> message.date
9 16092008

```

14.4 Les docstrings

Lorsqu'on souhaite distribuer un module, ou encore si on ne l'a pas utilisé depuis longtemps, il peut être utile de créer de la documentation pour expliquer ce que fait le module, et comment utiliser chaque fonction. Les commentaires (entre triple guillemets) situés en début de module et sous chaque fonction sont là pour cela, on les appelle *docstrings*. Ces *docstrings* permettent notamment de fournir de l'aide lorsqu'on invoque la commande `help()` :

```

1 >>> help(message)
2
3 Help on module message:
4
5 NAME
6     message - Module inutile qui affiche des messages :- )
7
8 FUNCTIONS
9     Bonjour(nom)
10         Affiche bonjour !
11

```

1. <https://docs.python.org/3.6/glossary.html>

```

12     Ciao(nom)
13         Affiche ciao !
14
15     Hello(nom)
16         Affiche hello !
17
18 DATA
19     date = 16092008
20
21 FILE
22     /home/atlas/fuchs/message.py

```

Vous voyez que Python a formaté automatiquement cette page d'aide, tout comme il est capable de le faire pour les modules internes à Python (e.g. `random`, `math`, etc), et ce grâce aux *docstrings*. Notez que l'on peut invoquer l'aide pour une seule fonction :

```

1 >>> help(message.Ciao)
2
3 Help on function Ciao in module message:
4
5 Ciao(nom)
6     Affiche ciao !

```

En résumé, les *docstrings* sont destinés aux utilisateurs du module, leur but est donc différent des commentaires qui eux sont destinés à celui qui lit le code (pour en comprendre les subtilités). Une bonne *docstring* de fonction doit contenir tout ce dont un utilisateur a besoin pour utiliser cette fonction. Une liste minimale (non exhaustive) serait :

- ce que fait la fonction,
- ce qu'elle prend en argument,
- ce qu'elle renvoie.

Pour en savoir plus sur les docstrings, vous pouvez consulter les recommandations officielles des développeurs de Python sur la page nommée PEP 257²

14.5 Modules et fonctions

La visibilité des fonctions au sein des modules suit des règles simples :

- Les fonctions dans un module peuvent s'appeler les unes les autres.
- Les fonctions dans un module peuvent appeler des fonctions situées dans un autre module s'il a été préalablement importé avec la commande `import`.

14.6 Module ou script ?

Vous avez remarqué que notre module `message` ne contient que des fonctions ou variables. Ainsi si nous le lançons comme un script classique, cela ne renverrait aucune sortie :

```

1 $ python ./message.py
2 $

```

2. <https://www.python.org/dev/peps/pep-0257/>

Cela vient bien sûr de l'absence de programme principal, c'est à dire, de ligne(s) de code située(s) en colonne 1 et que l'interpréteur exécute lorsqu'on lance le script.

Que se passe-t-il alors si on importe un script en tant que module alors qu'il contient un programme principal avec des lignes de code? Par exemple, avec le code suivant enregistré dans un fichier `message2.py` :

```

1 #!/usr/bin/env python3
2
3 def Bonjour(nom):
4     """Affiche bonjour !
5     """
6     return "bonjour " + nom
7
8 # programme principal
9 print(Bonjour("Joe"))

```

Si on l'importe dans l'interpréteur, on obtient :

```

1 >>> import message2
2 bonjour Joe

```

Ceci n'est pas le comportement attendu d'un module car on n'attend pas d'affichage particulier (par exemple un `import math` n'affiche rien dans l'interpréteur), et plus généralement on s'attend à avoir des fonctions accessibles mais pas spécialement des lignes de code exécutées.

Afin de pouvoir utiliser un code Python en tant que module ou en tant que script, il faudra utiliser la structure suivante :

```

1 #!/usr/bin/env python3
2
3 def Bonjour(nom):
4     """Affiche bonjour !
5     """
6     return "bonjour " + nom
7
8 if __name__ == '__main__':
9     print(Bonjour("Joe"))

```

La ligne `if __name__ == '__main__':` : située en colonne 1 indiquera à Python :

- si le programme `message2.py` a été exécuté en tant que script, le résultat du `if` sera alors `True` :

```

1 $ python3 ./message2.py
2 bonjour Joe
3 $

```

- si le programme `message2.py` a été importé en tant que module, le résultat du `if` sera alors `False` :

```

1 >>> import message2
2 >>>

```

dans ce dernier cas les fonctions seront accessibles mais le bloc du `if` ne sera pas exécuté. Ce comportement est possible grâce à la gestion de l'espace des noms par Python [AJOUTER ICI UN LIEN vers une section espace des noms].

Au delà de la commodité de pouvoir utiliser votre programme en tant que script ou en tant que module, cela présente l'avantage de bien voir où se situe le programme principal quand on lit le code (ainsi, plus besoin d'ajouter un commentaire `# programme principal` comme nous vous l'avions recommandé dans les chapitres 9 et 12 sur les fonctions). L'utilisation de la ligne `if __name__ == '__main__':` est ainsi une bonne pratique que nous vous recommandons !

14.7 Exercices

Conseil : pour cet exercice, écrivez un script dans un fichier, puis exécutez-le dans un *shell*.

14.7.1 Module ADN

Reprenez l'ensemble des fonctions qui gèrent le traitement de séquences nucléiques et incluez-les dans un module *adn.py*. Testez-les au fur et à mesure.

Chapitre 15

Expressions régulières et parsing

Le module `re`¹ vous permet d'utiliser des expressions régulières au sein de Python. Les expressions régulières sont aussi appelées en anglais *regular expressions* ou en plus court *regex* (dans la suite de ce chapitre, nous utiliserons souvent le mot *regex* pour désigner une expression régulière). Elles sont puissantes et donc incontournables en bioinformatique, spécialement lorsque vous souhaitez récupérer des informations dans des gros fichiers.

Cette action de recherche de données dans un fichier est appelée plus généralement *parsing* (qui signifie littéralement “analyse syntaxique”). Le *parsing* fait partie du travail quotidien du bioinformaticien, il est sans arrêt en train de “fouiller” dans des fichiers pour en extraire des informations d'intérêt comme par exemple récupérer les coordonnées 3D des atomes d'une protéines dans un fichier PDB ou alors extraire les gènes d'un fichier genbank.

Dans ce chapitre, nous ne ferons que quelques rappels sur les expressions régulières. Pour une documentation plus complète, référez-vous à la page d'aide des expressions régulières² sur le site officiel de Python.

15.1 Définition et syntaxe

Une expression régulière est une suite de caractères qui a pour but de décrire un fragment de texte. Cette suite de caractères est encore appelée *motif* (en anglais *pattern*), motif qui est constitué de deux types de caractères :

- Les caractères dits *normaux*.
- Les *métacaractères* ayant une signification particulière, par exemple le caractère `^` signifie début de ligne et non pas le caractère “chapeau” littéral.

Certains programmes Unix comme `egrep`, `sed` ou encore `awk` savent interpréter les expressions régulières. Tous ces programmes fonctionnent généralement selon le schéma suivant :

- Le programme lit un fichier ligne par ligne.
- Pour chaque ligne lue, si l'expression régulière passée en argument est retrouvée dans la ligne alors le programme effectue une action.

Par exemple, pour le programme `egrep` :

```
1 $ egrep "^DEF" herp_virus.gbk
2 DEFINITION Human herpesvirus 2, complete genome.
3 $
```

1. <https://docs.python.org/3/library/re.html>

2. <https://docs.python.org/3/library/re.html>

Ici, `egrep` affiche toutes les lignes du fichier `genbank` du virus de l'herpès (`herp_virus.gbk`) dans lesquelles la regex `^DEF` (c'est-à-dire le mot `DEF` en début de ligne) est retrouvée.

Il est intéressant ici de faire une petite digression sur le vocabulaire utilisé en anglais : en général, on utilise le verbe *to match* pour indiquer qu'une regex "a fonctionné". Bien qu'il n'y ait pas de traduction littérale en français, on peut utiliser les verbes "retrouver" ou "correspondre". Par exemple, pour traduire la phrase *The regex matches the line* on pourra utiliser "la regex est retrouvée dans la ligne" ou encore "la regex correspond dans la ligne".

Avant de voir comment Python gère les expressions régulières, voici quelques éléments de syntaxe des métacaractères :

`^` Début de chaîne de caractères ou de ligne.

Exemple : la regex `^ATG` est retrouvée dans la chaîne de caractères `ATGCGT` mais pas dans la chaîne `CCATGTT`.

`$` Fin de chaîne de caractères ou de ligne.

Exemple : la regex `ATG$` est retrouvée correspond dans la chaîne de caractères `TGCATG` mais pas dans la chaîne `CCATGTT`.

`.` N'importe quel caractère (mais un caractère quand même).

Exemple : la regex `A.G` est retrouvée dans `ATG`, `AtG`, `A4G`, mais aussi dans `A-G` ou dans `A G`.

`[ABC]` Le caractère `A` ou `B` ou `C` (un seul caractère).

Exemple : la regex `T[ABC]G` est retrouvée dans `TAG`, `TBG` ou `TCG`, mais pas à `TG`.

`[A-Z]` N'importe quelle lettre majuscule.

Exemple : la regex `C[A-Z]T` est retrouvée dans `CAT`, `CBT`, `CCT`...

`[a-z]` N'importe quelle lettre minuscule.

`[0-9]` N'importe quel chiffre.

`[A-Za-z0-9]` N'importe quel caractère alphanumérique.

`[^AB]` N'importe quel caractère sauf `A` et `B`.

Exemple : la regex `CG[^AB]T` est retrouvée dans `CG9T`, `CGCT`... mais pas dans `CGAT` ni dans `CGBT`.

`\` Caractère d'échappement (pour protéger certains caractères).

Exemple : la regex `\+` désigne le caractère `+` littéral. La regex `A\.G` est retrouvée dans `A.G` et non pas dans `A` suivi de n'importe quel caractère, suivi de `G`.

`*` 0 à n fois le caractère précédent ou l'expression entre parenthèses précédente.

Exemple : la regex `A(CG)*T` est retrouvée dans `AT`, `ACGT`, `ACGCGT`...

`+` 1 à n fois le caractère précédent ou l'expression entre parenthèses précédente.

Exemple : la regex `A(CG)+T` est retrouvée dans `ACGT`, `ACGCGT`... mais pas dans `AT`.

`?` 0 à 1 fois le caractère précédent ou l'expression entre parenthèses précédente.

Exemple : la regex `A(CG)?T` est retrouvée dans `AT` ou `ACGT`.

`{n}` n fois le caractère précédent ou l'expression entre parenthèses précédente.

Exemple : la regex `A(CG){2}T` est retrouvée dans `ACGCGT` mais pas dans `ACGT`, `ACGCGCGT` ou `ACGCG`.

`{n,m}` n à m fois le caractère précédent ou l'expression entre parenthèses précédente.

Exemple : la regex `A(C){2,4}T` est retrouvée dans `ACCT`, `ACCCT` et `ACCCCT` mais pas dans `ACT`, `ACCCCT` ou `ACCC`.

`{n,}` Au moins n fois le caractère précédent ou l'expression entre parenthèses précédente.

Exemple : la regex `A(C){2,}T` est retrouvée dans `ACCT`, `ACCCT` et `ACCCCT` mais pas à `ACT` ou `ACCC`.

`{,m}` Au plus m fois le caractère précédent ou l'expression entre parenthèses précédente.

Exemple : la regex `A(C){,2}T` est retrouvée dans `AT`, `ACT` et `ACCT` mais pas dans `ACCCT` ou `ACC`.

`(CG|TT)` Les chaînes de caractères `CG` ou `TT`.

Exemple : la regex `A(CG|TT)C` est retrouvée dans `ACGC` ou `ATTC`.

Enfin, il existe des caractères spéciaux qui sont bien commodes et qui peuvent être utilisés en tant que métacaractères :

`\d` remplace n'importe quel chiffre (d signifie *digit*), équivalent à `[0-9]`.

`\w` remplace n'importe quel caractère alphanumérique et le caractère souligné (*underscore*) (w signifie *word character*), équivalent à `[0-9A-Za-z_]`.

`\s` remplace n'importe quel "espace blanc" (*whitespace*) (s signifie *space*), équivalent à `[\t\n]` (c'est à dire espace, tabulation, nouvelle ligne), **attention** `\s` remplace également les caractères spéciaux `\r` et `\f` que nous ne développerons pas ici. `\s` est particulièrement pratique pour détecter une combinaison d'espace(s) et/ou de tabulation(s).

Comme vous le constatez, les métacaractères sont nombreux et leur signification est parfois difficile à maîtriser. Faites particulièrement attention aux métacaractères `.`, `+` et `*` qui combinés ensemble peuvent donner des résultats ambigus.

Il est important de savoir par ailleurs que les regex sont *avides* lorsqu'on utilise les métacaractères `+` et `*` lors d'une substitution. C'est à dire que la regex cherchera à "s'étendre" au maximum. Par exemple, si on utilise la regex `A+` pour faire une substitution dans la chaîne `TTTAAAAAAAAGC`, tous les `A` de cette chaîne (8 en tout) seront concernés, bien que `AA`, `AAA`, etc, "fonctionnent" également avec cette regex.

Nous vous conseillons de tester systématiquement vos expressions régulières sur des exemples simples. Pour vous aider nous vous recommandons des sites internet :

- <https://regexone.com/> qui propose un petit tutorial en ligne très bien fait.
- <https://regexr.com/> qui permet de visualiser tous les endroits où une regex est retrouvée dans un texte.
- <https://www.regular-expressions.info> vous donnera une documentation exhaustive sur les regex (il y a même une section sur Python).

N'hésitez pas à tester ces sites avant de vous lancer dans les exercices ou dans l'écriture de vos propres regex !

15.2 Module re et fonction search

15.2.1 La fonction search()

Dans le module `re`, la fonction incontournable est la fonction `search()`. Elle permet de rechercher un motif (*pattern* en anglais, c'est à dire une regex) au sein d'une chaîne de caractères avec une syntaxe de la forme `search(motif, chaine)`. Si `motif` est retrouvé dans `chaine`, Python renvoie un objet du type `SRE_Match`.

Attention, le motif (ou *pattern*) que vous utilisez comme premier argument de la fonction `search()` sera interprété *en tant que regex*. Ainsi, `^DEF` correspondra au mot `DEF` en début de chaîne, et non pas au caractère littéral `^` suivi du mot `DEF`.

Sans entrer dans les détails propres au langage orienté objet, si on utilise un objet du type `SRE_Match` dans un test, il sera considéré comme vrai. Regardez cet exemple dans lequel on va rechercher le motif `tigre` dans la chaîne de caractères `"girafe tigre singe"` :

```

1 >>> import re
2 >>> animaux = "girafe tigre singe"
3 >>> re.search('tigre', animaux)
4 <_sre.SRE_Match object at 0x7fefdaefe2a0>
5 >>> if re.search('tigre', animaux):
6 ...     print "OK"
7 ...
8 OK

```

15.2.2 Les fonctions match() et fullmatch()

Il existe aussi la fonction `match()` dans le module `re` qui fonctionne sur le modèle de `search()`. La différence est qu'elle renvoie un objet du type `SRE_Match` seulement lorsque l'expression régulière correspond au début de la chaîne de caractères (à partir du premier caractère).

```

1 >>> animaux = "girafe tigre singe"
2 >>> re.search('tigre', animaux)
3 <_sre.SRE_Match object at 0x7fefdaefe718>
4 >>> re.match('tigre', animaux)
5 >>>
6 >>> animaux = "tigre singe"
7 >>> re.match('tigre', animaux)
8 <_sre.SRE_Match object; span=(0, 5), match='tigre'>
9 >>>

```

Il existe également la fonction `fullmatch()` qui renvoie un objet du type `SRE_Match` si et seulement si l'expression régulière correspond **exactement** à la chaîne de caractères.

```

1 >>> animaux = "tigre "
2 >>> re.fullmatch('tigre', animaux)
3 >>> animaux = "tigre"
4 >>> re.fullmatch('tigre', animaux)
5 <_sre.SRE_Match object; span=(0, 5), match='tigre'>

```

De manière générale, nous vous recommandons plutôt l'usage de la fonction `search()`. Si vous souhaitez avoir une correspondance avec le début de la chaîne comme dans la fonction `match()`, vous pouvez toujours utiliser l'accroche de début de ligne `^`. Si vous voulez une correspondance exacte comme dans la fonction `fullmatch()`, vous pouvez utiliser les métacaractères `^` et `$`, par exemple `^tigre$`.

15.2.3 Compilation d'expressions régulières

Lorsqu'on a besoin de tester la même expression régulière sur plusieurs milliers de chaînes de caractères, il est pratique de préalablement compiler l'expression régulière à l'aide de la fonction `compile()` qui renvoie un objet de type `SRE_Pattern` :

```

1 >>> regex = re.compile("^tigre")
2 >>> regex
3 <_sre.SRE_Pattern object at 0x7fefdafd0df0>

```

On peut alors utiliser directement cet objet avec la méthode `search()` :

```

1 >>> animaux = "girafe tigre singe"
2 >>> regex.search(animaux)
3 >>> animaux = "tigre singe"
4 >>> regex.search(animaux)
5 <_sre.SRE_Match object at 0x7fefdaefe718>
6 >>> animaux = "singe tigre"
7 >>> regex.search(animaux)

```

15.2.4 Groupes

L'intérêt de l'objet de type `SRE_Match` renvoyé par Python lorsqu'une expression régulière trouve une correspondance dans une chaîne de caractères est de pouvoir ensuite récupérer certaines zones précises :

```

1 >>> regex = re.compile('[0-9]+\.[0-9]+')

```

Dans cet exemple, on recherche un nombre :

- qui débute par un ou plusieurs chiffres `[0-9]+`,
- suivi d'un point `\.` (le point a d'habitude une signification de métacaractère, donc il faut l'échapper avec `\` pour qu'il retrouve sa signification de point),
- et qui se termine encore par un ou plusieurs chiffres `[0-9]+`.

Les parenthèses dans l'expression régulière permettent de créer des groupes (`[0-9]+` deux fois) qui seront récupérés ultérieurement par la fonction `group()`.

```

1 >>> resultat = regex.search("pi vaut 3.14")
2 >>> resultat.group(0)
3 '3.14'
4 >>> resultat.group(1)
5 '3'
6 >>> resultat.group(2)
7 '14'
8 >>> resultat.start()
9 8
10 >>> resultat.end()
11 12

```

La totalité de la correspondance est donné par `group(0)`, le premier élément entre parenthèse est donné par `group(1)` et le second par `group(2)`.

Les fonctions `start()` et `end()` donnent respectivement la position de début et de fin de la zone qui correspond à l'expression régulière. Notez que la fonction `search()` ne renvoie que la première zone qui correspond à l'expression régulière, même s'il en existe plusieurs :

```

1 >>> resultat = regex.search("pi vaut 3.14 et e vaut 2.72")
2 >>> resultat.group(0)
3 '3.14'

```

15.2.5 La fonction findall()

Pour récupérer chaque zone, vous pouvez utiliser la fonction `findall()` qui renvoie une liste des éléments en correspondance.

```

1 >>> regex = re.compile('[0-9]+\.[0-9]+')
2 >>> resultat = regex.findall("pi vaut 3.14 et e vaut 2.72")
3 >>> resultat
4 ['3.14', '2.72']

```

L'utilisation des groupes entre parenthèse est également possible et ceux-ci sont automatiquement renvoyés sous la forme de tuples.

```

1 >>> regex = re.compile('([0-9]+)\.([0-9]+)')
2 >>> resultat = regex.findall("pi vaut 3.14 et e vaut 2.72")
3 >>> resultat
4 [( '3', '14'), ('2', '72')]

```

15.2.6 La fonction sub()

Enfin, la fonction `sub()` permet d'effectuer des remplacements assez puissants. Par défaut la fonction `sub(chaine1,chaine2)` remplace toutes les occurrences trouvées par l'expression régulière dans `chaine2` par `chaine1`. Si vous souhaitez ne remplacer que les n premières occurrences, utilisez l'argument `count=n` :

```

1 >>> regex = re.compile('[0-9]+\.[0-9]+')
2 >>> regex.sub('quelque chose',"pi vaut 3.14 et e vaut 2.72")
3 'pi vaut quelque chose et e vaut quelque chose'
4 >>> regex.sub('quelque chose',"pi vaut 3.14 et e vaut 2.72", count=1)
5 'pi vaut quelque chose et e vaut 2.72'

```

Encore plus puissant, il est possible d'utiliser dans le remplacement des groupes qui ont été "capturés" avec des parenthèses.

```

1 >>> regex = re.compile('([0-9]+)\.([0-9]+)')
2 >>> phrase = "pi vaut 3.14 et e vaut 2.72"
3 >>> regex.sub("approximativement \1",phrase)
4 'pi vaut approximativement 3 et e vaut vaut approximativement 2'
5 >>>
6 >>> regex.sub("approximativement \1 (puis .\2)",phrase)
7 'pi vaut approximativement 3 (puis .14) et e vaut approximativement 2 (puis .72)'

```

Si vous avez capturé des groupes, il suffit d'utiliser `\1`, `\2` (etc) pour utiliser les groupes correspondants dans la chaîne substituée. On pourra noter que la syntaxe générale pour récupérer des groupes dans les outils qui gèrent les regex est `\1`, `\2` (etc). Toutefois Python nous oblige à mettre un deuxième backslash car il y a ici deux niveaux : i) un premier niveau Python où on veut mettre un backslash littéral (donc `\\`), puis ii) un deuxième niveau regex dans lequel on veut retrouver `\1`. Si cela est confus, retenez seulement qu'il faut mettre un `\\` devant le numéro de groupe.

Nous espérons que vous êtes convaincus de la puissance du module `re` et des expressions régulières, alors plus de temps à perdre, à vos regex !

15.3 Exercices

15.3.1 Regex de base

Dans ces exercices, nous allons utiliser le fichier `genbank NC_001133.gbk`³ correspondant au chromosome I de la levure *Saccharomyces Cerevisiae* (souche S288C).

- Créer un script `def.py` qui recherche le mot `DEFINITION` dans le fichier `gbk`, et ce en début de ligne. Le script affichera cette ligne (équivalent de `egrep`)
- Ecrire un script `journal.py` qui affiche tous les journaux (mot-clé `JOURNAL`) dans lesquels ont été publié les travaux sur ce génome (bien sûr utiliser une regex!).

15.3.2 Nettoyeur d'espaces

Le fichier `cigale_fourmi.txt`⁴ contient le célèbre poème de Jean de la Fontaine. Malheureusement celui qui l'a recopié a parfois mis plusieurs espaces au lieu d'un seul entre les différents mots. Créer un script `cigale_fourmi.py` qui grâce aux regex et à la fonction `sub()` remplace les combinaisons de 2, 3, 4 (etc) espaces par un seul espace. Le nouveau texte "propre" sera enregistré dans un fichier `cigale_fourmi_propre.txt`.

15.3.3 Nettoyeur de doublons

Ecrire un script `ote_doublons.py` qui lit un fichier `breves_doublons.txt`⁵ et qui ôte tous les doublons de celui-ci à l'aide d'une regex. Le script affichera le nouveau texte à l'écran.

15.3.4 Compteur de gènes

- Télécharger un petit fichier `gbk` sur la genbank (vous pouvez réutiliser `NC_001133.gbk`).
- Ecrire un script `genes.py` qui affiche tous les mots `gene` dans le fichier `gbk`. Lancer ce script de la manière suivante afin de voir ce que celui-ci renvoie page par page :

```
1 $ python3 genes.py | less
```

Si on compte ces lignes est-ce une bonne méthode pour compter les gènes ? Comment selon vous peut-on récupérer une ligne unique par gène (bien observer le fichier `genbank`) ?

- Sur la base du script précédent, écrire un script `compte_genes.py` qui compte le nombre de gènes total et qui l'affiche à l'écran.
- Améliorer le script afin qu'il affiche en plus le nombre de gènes directs et le nombre de gènes complémentaires.

15.3.5 Recherche de séquence dans un fichier gbk

On veut créer un script `get_seq.py` qui extrait une séquence d'un fichier `genbank`.

- Télécharger un petit fichier `gbk` sur la genbank (vous pouvez réutiliser `NC_001133.gbk`).
- Afficher toutes les lignes qui commencent par au moins un espace suivi d'un chiffre. Cela correspond-il aux lignes contenant la séquence d'ADN ?

3. [data-files/NC_001133.gbk](#)

4. [data-files/cigale_fourmi.txt](#)

5. [data-files/breves_doublons.txt](#)

- Affiner la regex pour ne récupérer que les lignes contenant la séquence. Testez votre nouvelle regex sur d'autres fichiers genbank récupérés sur internet.
- Avec la regex trouvée précédemment, faites un script qui récupère la séquence (dans une chaîne de caractères) en récupérant les bonnes lignes et en substituant les chiffres, espaces et retours chariot par rien du tout. Le script fonctionnera comme suit :

```

1 $ python3 get_seq.py file.gbk
2 Le fichier file.gbk contient un génome de XXX bases.
3 $

```

où XXX correspond au nombre de paires de bases dans le génome. Comparez avec le nombre de bases indiqué dans le fichier genbank (ligne LOCUS).

15.3.6 Le défi du dé-htmliseur (+++)

Le format html permet d'afficher des pages web sur un navigateur. Il s'agit d'un langage à balise qui fonctionne avec des balises ouvrantes `<balise>` et des balises fermantes `</balise>`. Créer un script `dhtmliseur.py` qui lit le fichier `fichier_a_dhtmliser.html`⁶ et qui renvoie à l'écran tout le texte de ce fichier contenu en dehors des balises html. Nous vous conseillons d'ouvrir le fichier html dans un éditeur de texte et de bien l'observer. N'hésitez pas à vous aider du site <https://regexr.com/>.

15.3.7 Remplacements puissants (+++)

A partir du fichier contenant le protéome humain⁷, on souhaite faire une liste de toutes les protéines et les indexer avec un numéro arbitraire. Nous allons pour cela utiliser les regex et la substitution. Ecrire un script `liste_proteome.py` qui procède comme suit :

- Lire le fichier de protéome, et afficher uniquement les lignes de commentaires du fichier multi-fasta à l'aide d'une regex (lignes commençant par `>`).
- Développez votre regex afin de récupérer le numéro d'accèsion dans un groupe (par exemple dans la ligne `>sp|095139|NDUB6_HUMAN NADH dehydrogenase [...]` le numéro d'accèsion est `095139`, il se situe entre le premier et le deuxième symbole `|` (symbole *pipe*)).
- Au lieu d'afficher les lignes complète, afficher seulement le numéro d'accèsion, à l'aide de la fonction `sub()`.
- Au final, on souhaite la sortie suivante (pensez à l'opérateur de formatage qui ajoute des 0 à la place d'espaces, par exemple `{:03d}`) :

```

1 protein 00001 095139
2 protein 00002 075438
3 [...]
4 protein 20372 Q9UKP6
5 protein 20373 Q96HZ7

```

6. [data-files/fichier_a_dhtmliser.html](#)

7. [data-files/human-proteome.fasta](#)

Chapitre 16

Autres modules d'intérêt

en construction

Chapitre 17

Modules d'intérêt en bioinformatique

Nous allons découvrir dans cette section quelques modules très importants en bioinformatique. Le premier `numpy` permet notamment de manipuler des vecteurs et des matrices en Python. Le module `biopython` permet de travailler sur des données biologiques, comme des séquences (nucléiques et protéiques) ou des structures (fichiers PDB). Le module `matplotlib` permet de dessiner des graphiques depuis Python.

Ces modules ne sont pas fournis avec la distribution Python de base (contrairement à tous les autres modules vus précédemment). Nous ne nous étendrons pas sur la manière de les installer. Consultez pour cela la documentation sur les sites internet des modules en question. Sachez cependant que ces modules existent dans la plupart des distributions Linux récentes.

Dans ce chapitre, nous vous montrerons quelques exemples d'utilisation de ces modules pour vous convaincre de leur pertinence.

17.1 Module numpy

Le module `numpy`¹ est incontournable en bioinformatique. Il permet d'effectuer des calculs sur des vecteurs ou des matrices, élément par élément, via un nouveau type d'objet appelé *array*. Ce module contient des fonctions de base pour faire de l'algèbre linéaire, des transformées de Fourier ou encore des tirages de nombre aléatoire plus sophistiqués qu'avec le module `random`. Vous pouvez télécharger `numpy` à cette adresse². Notez qu'il existe un autre module `scipy`³ que nous n'aborderons pas dans ce cours. `scipy` est lui même basé sur `numpy`, mais il en étend considérablement les possibilités de ce dernier (par exemple : statistiques, optimisation, intégration numérique, traitement du signal, traitement d'image, algorithmes génétiques, etc).

On charge le module `numpy` avec la commande

```
1 >>> import numpy
```

On peut également définir un raccourci pour `numpy` :

```
1 >>> import numpy as np
```

1. <http://numpy.scipy.org/>
2. <http://www.scipy.org/Download>
3. <http://www.scipy.org/>

17.1.1 Objets de type array

Les objets de type *array* correspondent à des tableaux à une ou plusieurs dimensions et permettent d'effectuer du calcul vectoriel. La fonction `array()` permet la conversion d'un objet séquentiel (type liste ou tuple) en un objet de type *array*. Voici un exemple simple de conversion d'une liste à une dimension en objet *array* :

```

1 >>> import numpy as np
2 >>> a = [1,2,3]
3 >>> np.array(a)
4 array([1, 2, 3])
5 >>> b = np.array(a)
6 >>> type(b)
7 <type 'numpy.ndarray'>
8 >>> b
9 array([1, 2, 3])

```

Nous avons converti la liste `a` en *array*, mais cela aurait donné le même résultat si on avait converti le tuple `(1,2,3)`. Par ailleurs, vous voyez que lorsqu'on demande à Python le contenu d'un objet *array*, les symboles `[` et `]` sont utilisés pour le distinguer d'une liste (délimitée par les caractères `[` et `]`) ou d'un tuple (délimité par les caractères `(` et `)`).

Notez qu'un objet *array* ne peut contenir que des valeurs numériques. Vous ne pouvez pas, par exemple, convertir une liste contenant des chaînes de caractères en objet de type *array*.

Contrairement à la fonction `range()`, la fonction `arange()` permet de construire un *array* à une dimension de manière simple.

```

1 >>> np.arange(10)
2 array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
3 >>> np.arange(10.0)
4 array([ 0.,  1.,  2.,  3.,  4.,  5.,  6.,  7.,  8.,  9.])
5 >>> np.arange(10,0,-1)
6 array([10, 9, 8, 7, 6, 5, 4, 3, 2, 1])

```

Un des avantages de la fonction `arange()` est qu'elle permet de générer des objets *array* qui contiennent des entiers ou de réels selon l'argument qu'on lui passe.

La différence fondamentale entre un objet *array* à une dimension et une liste (ou un tuple) est que celui-ci est considéré comme un **vecteur**. Par conséquent on peut effectuer des opérations **élément par élément** dessus, ce qui est bien commode lorsqu'on analyse de grandes quantités de données. Regardez ces exemples :

```

1 >>> v = np.arange(4)
2 >>> v
3 array([0, 1, 2, 3])
4 >>> v + 1
5 array([1, 2, 3, 4])
6 >>> v + 0.1
7 array([ 0.1,  1.1,  2.1,  3.1])
8 >>> v * 2
9 array([0, 2, 4, 6])
10 >>> v * v
11 array([0, 1, 4, 9])

```

Notez bien sur le dernier exemple de multiplication que l'*array* final correspond à la multiplication **élément par élément** des deux *array* initiaux. Avec les listes, ces opérations n'auraient été possibles qu'en utilisant des boucles! Nous vous encourageons donc à utiliser dorénavant les objets *array* lorsque vous aurez besoin de faire des opérations élément par élément.

17.1.2 array et dimensions

Il est aussi possible de construire des objets *array* à deux dimensions, il suffit de passer en argument une liste de listes à la fonction `array()` :

```
1 >>> w = np.array([[1,2], [3,4], [5,6]])
2 >>> w
3 array([[1, 2],
4        [3, 4],
5        [5, 6]])
```

On peut aussi créer des tableaux à trois dimensions en passant à la fonction `array()` une liste de listes de listes :

```
1 >>> x = np.array([[[1,2], [2,3]], [[4,5], [5,6]]])
2 >>> x
3 array([[[1, 2],
4         [2, 3]],
5        [[4, 5],
6         [5, 6]]])
7
```

La fonction `array()` peut créer des tableaux à n'importe quel nombre de dimensions. Toutefois ça devient vite compliqué lorsqu'on dépasse trois dimensions. Retenez qu'un objet *array* à une dimension peut être considéré comme un **vecteur** et un *array* à deux dimensions comme une **matrice**.

La méthode `.shape` renvoie les dimensions d'un *array* alors que `.size` renvoie le nombre d'éléments contenus dans l'*array*.

```
1 >>> v.shape
2 (4,)
3 >>> v.size
4 4
5 >>> w.shape
6 (3, 2)
7 >>> w.size
8 6
9 >>> x.shape
10 (2, 2, 2)
11 >>> x.size
12 8
```

Et la fonction `reshape()` permet de modifier les dimensions d'un *array* :

```
1 >>> a = np.arange(0, 6)
2 >>> a
3 array([0, 1, 2, 3, 4, 5])
4 >>> a.shape
5 (6,)
```

```

6 >>> b = t.reshape((2, 3))
7 >>> b
8 array([[0, 1, 2],
9        [3, 4, 5]])
10 >>> b.shape
11 (2, 3)

```

Notez que `a.reshape((2, 3))` n'est pas la même chose que `a.reshape((3, 2))` :

```

1 >>> c = a.reshape((3, 2))
2 >>> c
3 array([[0, 1],
4        [2, 3],
5        [4, 5]])
6 >>> c.shape
7 (3, 2)

```

La fonction `reshape()` attend que les nouvelles dimensions soient **compatibles** avec la dimension initiale de l'objet *array*, c'est-à-dire que le nombre d'éléments contenus dans les différents *array* soit le même. Dans nos exemples précédents, $6 = 2 \times 3 = 3 \times 2$.

Si les nouvelles dimensions ne sont pas compatibles avec les dimensions initiales, `reshape()` génère une erreur.

```

1 >>> a = np.arange(0, 6)
2 >>> a.shape
3 (6,)
4 >>> d = a.reshape((3, 4))
5 Traceback (most recent call last):
6 File "<stdin>", line 1, in <module>
7 ValueError: cannot reshape array of size 6 into shape (3,4)

```

La fonction `resize()` par contre ne déclenche pas d'erreur dans une telle situation et ajoute des 0 ou coupe la liste initiale jusqu'à ce que le nouvel *array* soit rempli.

```

1 >>> a = np.arange(0, 6)
2 >>> a.shape
3 (6,)
4 >>> a.resize((3,3))
5 >>> a.shape
6 (3, 3)
7 >>> a
8 array([[0, 1, 2],
9        [3, 4, 5],
10       [0, 0, 0]])
11
12
13 >>> b = np.arange(0, 10)
14 >>> b.shape
15 (10,)
16 >>> b.resize((2,3))
17 >>> b.shape
18 (2, 3)
19 >>> b

```

```

20 array([[0, 1, 2],
21        [3, 4, 5]])

```

À noter qu'il existe aussi la fonction `np.resize()` qui dans le cas d'un nouvel *array* plus grand que l'*array* initial, va répéter l'*array* initial :

```

1 >>> a = np.arange(0, 6)
2 >>> a.shape
3 (6,)
4 >>> c = np.resize(a, (3, 5))
5 >>> c.shape
6 (3, 5)
7 >>> c
8 array([[0, 1, 2, 3, 4],
9        [5, 0, 1, 2, 3],
10       [4, 5, 0, 1, 2]])

```

17.1.3 Indices

Pour récupérer un ou plusieurs élément(s) d'un objet *array*, vous pouvez utiliser l'indilage ou les tranchage, de la même manière que pour les listes.

```

1 >>> a = np.arange(10)
2 >>> a
3 array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
4 >>> a[5:]
5 array([5, 6, 7, 8, 9])
6 >>> a[::2]
7 array([0, 2, 4, 6, 8])
8 >>> a[1]
9 1

```

Dans le cas d'un objet *array* à deux dimensions, vous pouvez récupérer une ligne (d'indice *m*), une colonne (d'indice *n*) ou bien un seul élément.

```

1 >>> a = np.array([[1,2],[3,4]])
2 >>> a
3 array([[1, 2],
4        [3, 4]])
5 >>> a[:,0]
6 array([1, 3])
7 >>> a[0,:]
8 array([1, 2])
9 >>> a[1,1]
10 4

```

La syntaxe `a[m,:]` récupère la ligne *m-1*, et `a[:,n]` récupère la colonne *n-1*. Les tranches sont évidemment aussi utilisables sur un tableau à deux dimensions.

17.1.4 Construction automatique de matrices

Il peut être parfois pénible de construire une matrice (*array* à deux dimensions) à l'aide d'une liste de listes. Le module `numpy` contient quelques fonctions commodes pour construire des matrices à partir de

rien. Les fonctions `zeros()` et `ones()` permettent de construire des objets `array` contenant des 0 ou de 1, respectivement. Il suffit de leur passer un tuple indiquant les dimensions voulues.

```

1 >>> np.zeros((3,3))
2 array([[ 0.,  0.,  0.],
3        [ 0.,  0.,  0.],
4        [ 0.,  0.,  0.]])
5 >>> np.zeros((3,3), int)
6 array([[0, 0, 0],
7        [0, 0, 0],
8        [0, 0, 0]])
9 >>> np.ones((3,3))
10 array([[ 1.,  1.,  1.],
11         [ 1.,  1.,  1.],
12         [ 1.,  1.,  1.]])

```

Par défaut, les fonctions `zeros()` et `ones()` génèrent des réels, mais vous pouvez demander des entiers en passant l'option `int` en second argument.

Enfin, si voulez construire une matrice avec autre chose que des 0 ou des 1, vous avez à votre disposition la fonction `full()` :

```

1 >>> np.full((2,3), 7, int)
2 array([[7, 7, 7],
3        [7, 7, 7]])
4 >>> np.full((2,3), 7, float)
5 array([[ 7.,  7.,  7.],
6        [ 7.,  7.,  7.]])

```

Nous construisons ainsi une matrice constituée de 2 lignes et 3 colonnes et qui ne contient que le chiffre 7. Des entiers (`int`) dans le premier cas et des réels (`float`) dans le second.

17.1.5 Un peu d'algèbre linéaire

Après avoir manipulé les objets `array` comme des vecteurs et des matrices, voici quelques fonctions pour faire de l'algèbre linéaire.

La fonction `transpose()` renvoie la transposée d'un `array`. Par exemple pour une matrice :

```

1 >>> a
2 array([[1, 2, 3],
3        [4, 5, 6],
4        [7, 8, 9]])
5 >>> np.transpose(a)
6 array([[1, 4, 7],
7        [2, 5, 8],
8        [3, 6, 9]])

```

La fonction `dot()` vous permet de faire une multiplication de matrices.

```

1 >>> a = np.resize(numpy.arange(4), (2,2))
2 >>> a
3 array([[0, 1],
4        [2, 3]])
5 >>> np.dot(a,a)

```



```

6 array([[ 2,  3],
7        [ 6, 11]])
8 >>> a * a
9 array([[0,  1],
10        [4,  9]])

```

Notez bien que `dot(a,a)` renvoie le **produit matriciel** entre deux matrices, alors que `a * a` renvoie le produit **élément par élément**.

Remarque : Dans numpy, il existe également des objets de type *matrix* pour lesquels les multiplications de matrices sont différents, mais nous ne les aborderons pas ici.

Pour toutes les opérations suivantes, il faudra utiliser des fonctions dans le sous-module `numpy.linalg`. La fonction `inv()` renvoie l'inverse d'une matrice carrée, `det()` son déterminant, `eig()` ses vecteurs et valeurs propres.

```

1 >>> a
2 array([[0,  1],
3        [2,  3]])
4 >>> np.linalg.inv(a)
5 array([[-1.5,  0.5],
6        [ 1. ,  0. ]])
7 >>> np.linalg.det(a)
8 -2.0
9 >>> np.linalg.eig(a)
10 (array([-0.56155281,  3.56155281]), array([[-0.87192821, -0.27032301],
11        [ 0.48963374, -0.96276969]]))
12 >>> np.linalg.eig(a)[0]
13 array([-0.56155281,  3.56155281])
14 >>> np.linalg.eig(a)[1]
15 array([[-0.87192821, -0.27032301],
16        [ 0.48963374, -0.96276969]])

```

Notez que la fonction `eig()` renvoie un tuple dont le premier élément correspond aux valeurs propres et le second élément aux vecteurs propres.

17.1.6 Un peu de transformée de Fourier

La transformée de Fourier est très utilisée pour l'analyse de signaux, notamment lorsqu'on souhaite extraire des périodicités au sein d'un signal bruité. Le module `numpy` possède la fonction `fft()` (dans le sous-module `fft`) permettant de calculer des transformées de Fourier.

Voici un petit exemple sur la fonction cosinus de laquelle on souhaite extraire la période à l'aide de la fonction `fft()` :

```

1 # 1) on définit la fonction y = cos(x)
2 import numpy as np
3 debut = -2 * np.pi
4 fin = 2 * np.pi
5 pas = 0.1
6 x = np.arange(debut,fin,pas)
7 y = np.cos(x)
8
9 # 2) on calcule la transformée de Fourier (TF) de la fonction cosinus

```

```

10 TF = np.fft.fft(y)
11 ABSTF = np.abs(TF)
12 # abscisse du spectre en radian-1
13 pas_xABSTF = 1/(fin-debut)
14 x_ABSTF = np.arange(0,pas_xABSTF * len(ABSTF),pas_xABSTF)

```

Plusieurs commentaires sur cet exemple :

- Vous constatez que `numpy` redéfinit certaines fonctions ou constantes mathématiques de base, comme `pi` (nombre π), `cos()` (fonction cosinus) ou `abs()` (valeur absolue, ou module d'un complexe). Ceci est bien pratique car nous n'avons pas à appeler ces fonctions ou constantes depuis le module `math`, le code en est ainsi plus lisible.
- Dans la partie 1, on définit le vecteur `x` représentant un angle allant de -2π à 2π radians par pas de 0,1 et le vecteur `y` comme le cosinus de `x`.
- Dans la partie 2, on calcule la transformée de Fourier avec la fonction `fft()` qui renvoie un vecteur (objet `array` à une dimension) de nombres complexes. Eh oui, le module `numpy` gère aussi les nombres complexes! On extrait ensuite le module du résultat précédent avec la fonction `abs()`.
- La variable `x_ABSTFL` représente l'abscisse du spectre (en radian^{-1}).
- La variable `ABSTF` contient le spectre lui même. L'analyse de ce dernier nous donne un pic à 0,15 radian^{-1} , ce qui correspond bien à 2π (plutôt bon signe de retrouver ce résultat). Le graphe de ce spectre est présenté dans la partie dédiée à `matplotlib`.

Notez que tout au long de cette partie, nous avons toujours utilisé la syntaxe `np.fonction()` pour bien vous montrer quelles étaient les fonctions propres à `numpy` (dont le nom est ici raccourcis par `np`).

17.2 Module biopython

Le module `biopython`⁴ propose de nombreuses fonctionnalités très utiles en bioinformatique. Le tutoriel⁵ est particulièrement bien fait, n'hésitez pas à le consulter.

Voici quelques exemples d'utilisation.

Définition d'une séquence :

```

1 >>> import Bio
2 >>> from Bio.Seq import Seq
3 >>> from Bio.Alphabet import IUPAC
4 >>> ADN = Seq("ATATCGGCTATAGCATGCA", IUPAC.unambiguous_dna)
5 >>> ADN
6 Seq('ATATCGGCTATAGCATGCA', IUPACUnambiguousDNA())

```

L'expression `IUPAC.unambiguous_dna` signifie que la séquence entrée est bien une séquence d'ADN.

Obtention de la séquence complémentaire et complémentaire inverse :

```

1 >>> ADN.complement()
2 Seq('TATAGCCGATATCGTACGT', IUPACUnambiguousDNA())
3 >>> ADN.reverse_complement()
4 Seq('TGCATGCTATAGCCGATAT', IUPACUnambiguousDNA())

```

Traduction en séquence protéique :

-
4. <http://biopython.org/>
 5. <http://biopython.org/DIST/docs/tutorial/Tutorial.html>

```

1 >>> ADN.translate()
2 Seq('ISAIAC', IUPACProtein())

```

17.3 Module matplotlib

Le module matplotlib⁶ permet de générer des graphiques depuis Python. Il est l'outil complémentaire de `numpy` et `scipy` lorsqu'on veut faire de l'analyse de données.

17.3.1 Représentation sous forme de points

Dans cet exemple, nous considérons l'évolution de la concentration d'un produit dans le sang (exprimé en mg/L) en fonction du temps (exprimé en heure).

Nous avons les résultats suivants :

Temps (h)	Concentration (mg/L)
1	3.5
2	5.8
3	9.1
4	11.8
6	17.5
7	21.3
9	26.8

Nous allons maintenant représenter l'évolution de la concentration en fonction du temps :

```

1 import matplotlib.pyplot as plt
2
3 temps = [1, 2, 3, 4, 6, 7, 9]
4 concentration = [5.5, 7.2, 11.8, 13.6, 19.1, 21.7, 29.4]
5
6 plt.scatter(temps, concentration, marker='o', color = 'blue')
7
8 plt.xlabel("Temps (h)")
9 plt.ylabel("Concentration (mg/L)")
10 plt.title("Concentration de produit en fonction du temps")
11 plt.show()

```

Vous devriez obtenir une fenêtre graphique **interactive** qui vous permet de manipuler le graphe (se déplacer, zoomer, enregistrer comme image, etc.) comme la figure 17.1.

Revenons maintenant sur le code. Tout d'abord, vous voyez qu'on importe le sous-module `pyplot` du module `matplotlib` et qu'on lui donne un nom court `plt` pour l'utiliser plus rapidement ensuite.

- Dans un premier temps, nous définissons les variables `temps` et `concentration` comme des listes. Les deux listes doivent avoir la même longueur (7 éléments dans le cas présent).
- La fonction `scatter()` permet de représenter des points sous forme de nuage de points. Les deux premiers arguments correspondent aux valeurs en abscisse et en ordonnée des points

6. <https://matplotlib.org/>

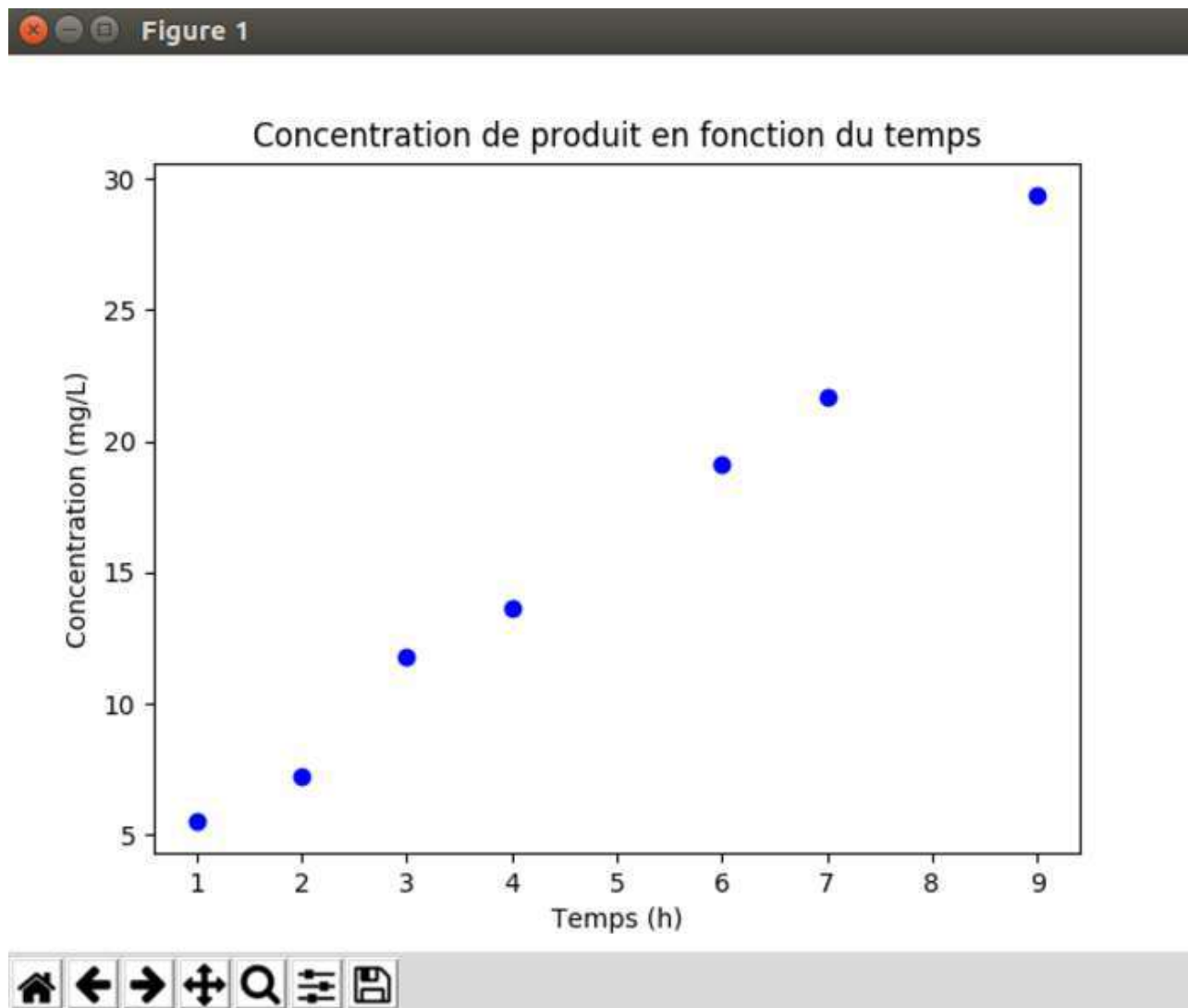


FIGURE 17.1 – Fenêtre interactive de matplotlib

expérimentaux, fournis sous forme de listes. Nous avons ensuite des arguments facultatifs comme le symbole (`marker`) et la couleur (`color`).

- Les fonctions `xlabel()` et `ylabel()` sont utiles pour donner un nom aux axes.
- La fonction `title()` permet de définir le titre du graphique.
- Enfin, la fonction `show()` affiche le graphique généré à l'écran.

17.3.2 Représentation sous forme de courbe

Nous savons par ailleurs que l'évolution de la concentration du produit en fonction du temps peut-être modélisée par la fonction $f(x) = 2 + 3 \times x$. Nous allons représenter cette fonction avec les points expérimentaux et sauvegarder le graphique obtenu sous forme d'une image :

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 temps = [1, 2, 3, 4, 6, 7, 9]
5 concentration = [5.5, 7.2, 11.8, 13.6, 19.1, 21.7, 29.4]
6
7 plt.scatter(temps, concentration, marker='o', color = 'blue')
8
9 plt.xlabel("Temps (h)")
10 plt.ylabel("Concentration (mg/L)")
11 plt.title("Concentration de produit en fonction du temps")
12
13 x = np.linspace( min(temps), max(temps), 50)
14 y = 2 + 3 * x
15 plt.plot(x, y, color='green', ls="--")
16 plt.grid()
17
18 plt.savefig('concentration_vs_temps.png', bbox_inches='tight', dpi=200)

```

Le résultat est représenté sur la figure 17.2.

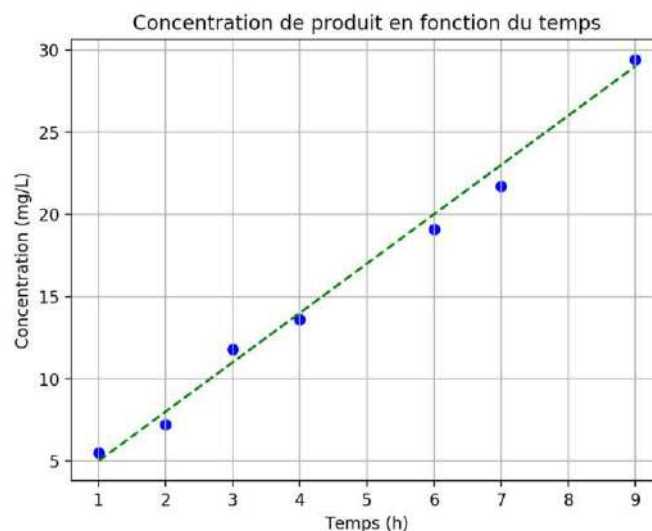


FIGURE 17.2 – Concentration du produit en fonction du temps

Nous avons ajouté les étapes suivantes :

- Nous chargeons le module `numpy` sous le nom `np`.
- Nous créons la variable `x` avec la fonction `linspace()` du module `numpy` qui crée une liste de valeurs régulièrement espacées entre deux bornes, ici le minimum (`min(temps)`) et le maximum (`max(temps)`) de la variable `temps`. Dans notre exemple, nous générons une liste de 50 valeurs. La variable `x` ainsi créée est du type `array`.
- Nous construisons ensuite la variable `y` à partir de la formule modélisant l'évolution de la concentration en fonction du temps. Cette manipulation n'est possible que parce que `x` est du type `array`. Cela ne fonctionnerait pas avec une liste classique.
- La fonction `plot()` permet de construire une courbe à partir des coordonnées en abscisse et en ordonnées des points à représenter. Nous avons ensuite des arguments facultatifs comme le style de la ligne (`ls`) et la couleur (`color`).
- La fonction `grid()` affiche une grille.
- Enfin, la fonction `savefig()` permet d'enregistrer le graphique produit sous la forme d'une image au format png. Des arguments optionnels définissent la manière de générer les marges autour du graphique (`bbox_inches`) et la résolution de l'image (`dpi`).

17.3.3 Représentation sous forme de barres

Nous souhaitons maintenant représenter graphiquement la distribution des différentes bases dans une séquence nucléique.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # 1 : calcul de la distribution des bases
5 sequence = "ACGATCATAGCGAGCTACGTAGAA"
6 bases = ["A", "C", "G", "T"]
7 distribution = []
8 for base in bases:
9     distribution.append(sequence.count(base))
10
11 # 2 : position des barres
12 x = np.arange(len(bases))
13
14 # 3 : représentation des barres
15 # et de l'axe des abscisses
16 plt.bar(x, distribution)
17 plt.xticks(x, bases)
18
19 # 4 : légendes des axes et titre du graphique
20 plt.xlabel("Bases")
21 plt.ylabel("Nombre")
22 plt.title("Distribution des bases\n dans la séquence {}".format(sequence))
23
24 # 5 : sauvegarde du graphique
25 plt.savefig('distribution_bases.png', bbox_inches='tight', dpi=200)

```

On obtient alors le graphique de la figure [17.3](#).

Prenons le temps d'examiner les différentes étapes du script précédent :

1. Nous définissons les variables (`sequence` et `bases`) que nous allons utiliser. Nous calculons ensuite la distribution des différentes bases dans la séquence. Nous utilisons pour cela la fonction

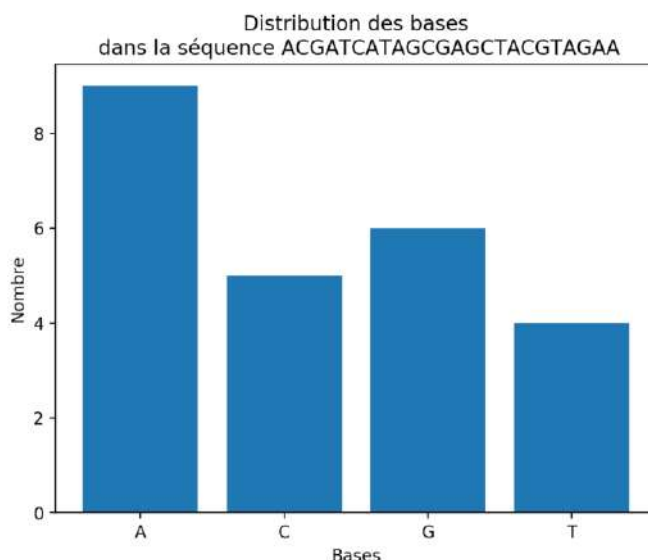


FIGURE 17.3 – Distribution des bases

- `count()` qui renvoie le nombre de fois qu'une sous-chaîne de caractères se trouve dans une autre.
2. Nous définissons la positions en abscisse des barres. Dans notre exemple, la variable `x` vaut `array([0, 1, 2, 3])`.
 3. Nous construisons la représentation graphique sous forme de barres avec la fonction `bar()` qui prend en argument la position des barres (`x`) et leurs hauteurs (`distribution`). La fonction `xticks()` redéfinit les étiquettes (c'est-à-dire le nom des bases) correspondantes aux différents positions des barres.
 4. Nous définissons les légendes des axes et le titre du graphique. Nous insérons un retour à la ligne `\n` dans le titre pour que cela soit plus esthétique.
 5. Enfin, nous enregistrons le graphique généré au format `png`.

Voilà, nous espérons que ce petit exemple vous aura convaincu de l'utilité du module `matplotlib`. Sachez qu'il peut faire bien plus, par exemple générer des histogrammes ou toutes sortes de graphiques utiles en analyse de données.

17.4 Exercices

Cet exercice présente le calcul de la distance entre deux atomes carbones alpha consécutifs de la barstar. Il demande quelques notions d'Unix et nécessite le module `numpy` de Python.

17.4.1 Extraction des coordonnées atomiques

Téléchargez le fichier `1bta.pdb` qui correspond à la structure tridimensionnelle de la protéine barstar⁷ sur le site de la PDB (lien direct vers le fichier⁸).

Voici le code pour extraire les coordonnées atomiques des carbones alpha de la barstar :

```
1 with open("1bta.pdb", "r") as fichier_pdb, open("1bta_ca.txt", "w") as fichier_CA:
2     for ligne in fichier_pdb:
```

7. <http://www.rcsb.org/pdb/explore.do?structureId=1BTA>

8. <https://files.rcsb.org/download/1BTA.pdb>

```

3     champs = ligne.split()
4     if champs[0] == "ATOM" and champs[2] == "CA":
5         fichier_CA.write("{} {} {} ".format(champs[6], champs[7], champs[8]))

```

Notez la structure de la forme

```

1 with open("1bta.pdb", "r") as fichier_pdb, open("1bta_CA.txt", "w") as fichier_CA:

```

qui permet d'ouvrir deux fichiers simultanément. Ici, le fichier `1bta.pdb` est ouvert en lecture (`r`) et le fichier `1bta_ca.txt` est ouvert en écriture (`w`).

Pour chaque ligne du fichier `pdb`, nous allons séparer le contenu de la ligne en champs et lorsque le premier champs vaut `ATOM` et le troisième vaut `CA`, alors nous écrivons les coordonnées de l'atome (7e, 8e et 9e champs) dans le fichier `1bta_ca.txt`. Les coordonnées sont toutes enregistrées sur une seule ligne, les unes après les autres.

Le même résultat peut être obtenu avec une seule commande Unix. Il faut par contre s'assurer que le séparateur décimal est bien le point (au lieu de la virgule, par défaut sur les systèmes d'exploitation français), et au préalable redéfinir la variable `LC_NUMERIC` en bash :

```

1 export LC_NUMERIC=C

```

La commande à utiliser est ensuite :

```

1 awk '$1=="ATOM" && $3=="CA" {printf "%.3f %.3f %.3f ", $7, $8, $9}' 1bta.pdb > 1bta_ca.txt

```

17.4.2 Lecture des coordonnées

Ouvrez le fichier `1bta_ca.txt` avec Python et créez une liste contenant toutes les coordonnées sous forme de réels avec les fonctions `split()` et `float()`.

Affichez à l'écran le nombre de total de coordonnées.

17.4.3 Utilisation de numpy

En ouvrant dans un éditeur de texte le fichier `1bta.pdb`, trouvez le nombre d'acides aminés qui constituent la barstar.

Avec la fonction `array()` du module `numpy`, convertissez la liste de coordonnées en `array`. Avec la fonction `reshape()` de `numpy`, construisez ensuite une matrice à deux dimensions contenant les coordonnées des carbones alpha de la barstar. Affichez les dimensions de cette matrice.

17.4.4 Calcul de la distance

Créez maintenant une matrice qui contient les coordonnées des `n-1` premiers carbones alpha et une autre qui contient les coordonnées des `n-1` derniers carbones alpha. Affichez les dimensions des matrices pour vérification.

En utilisant les opérateurs mathématiques habituels (`-`, `+`, `**2`) et les fonctions `sqrt()` et `sum()` du module `numpy`, calculez la distance entre les atomes `n` et `n+1`.

Pour chaque atome, affichez le numéro de l'atome et la distances entre carbones alpha consécutifs avec un chiffres après la virgule. Repérez la valeur surprenante.

Chapitre 18

Jupyter et ses notebooks

Les notebooks Jupyter sont des cahiers électroniques qui, dans le même document, peuvent rassembler du texte, des images et du code informatique exécutable. Ils sont manipulables interactivement dans un navigateur web.

Initialement développé pour les langages de programmation Julia, Python et R (d'où son nom), les notebooks Jupyter supportent près de 40 langages différents.

La cellule est l'élément de base d'un notebook Jupyter. Elle peut contenir du texte formaté au format Markdown ou du code informatique qui pourra être exécuté.

Voici un exemple de notebook Jupyter (figure 18.1) :

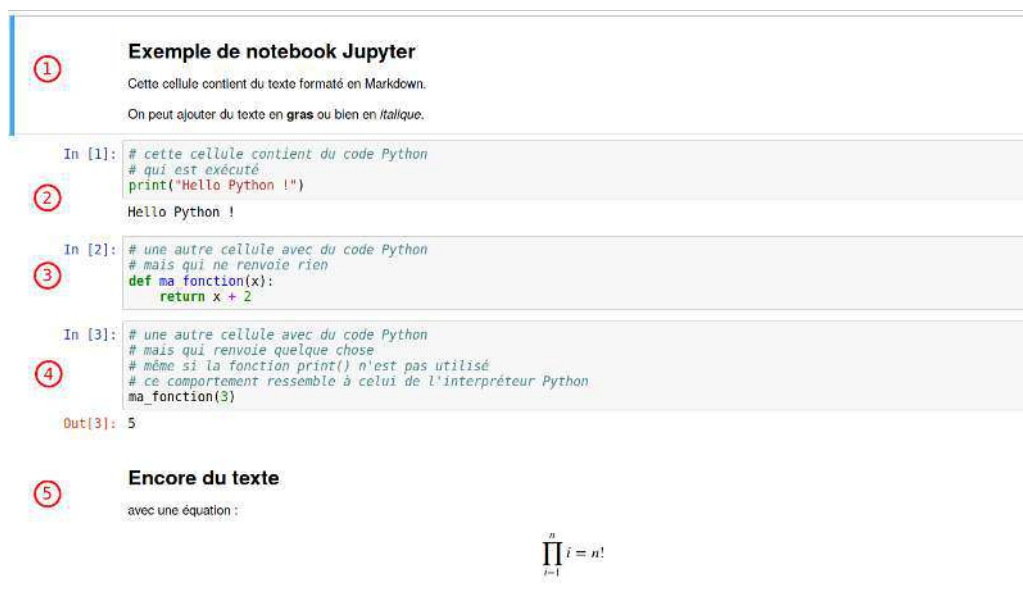


FIGURE 18.1 – Exemple de notebook Jupyter. Les chiffres dans les ronds indiquent les différentes cellules

Ce notebook est constitué de 5 cellules : deux avec du texte en Markdown (la première et la dernière) et 3 avec du code Python (notées avec In []).

18.1 Lancement de Jupyter et création d'un notebook

Jupyter est déjà installé avec la distribution Anaconda, pour le lancer, utilisez la commande :

```
1 $ jupyter-notebook
```

Une nouvelle page devrait s'ouvrir dans votre navigateur web et ressembler à la figure 18.2.

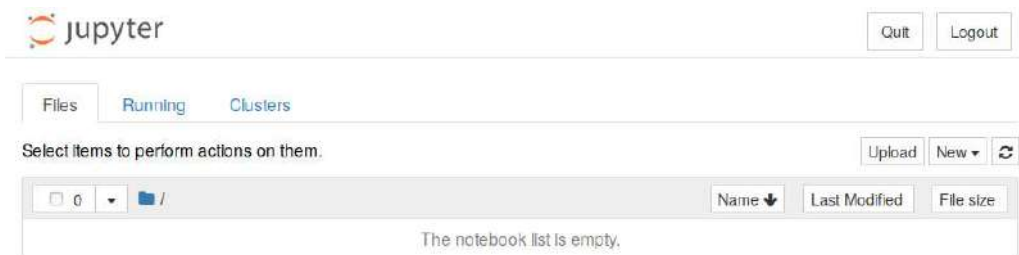


FIGURE 18.2 – Interface de Jupyter

Cette interface liste les notebooks Jupyter existants (pour le moment aucun).

Pour créer un notebook, cliquez sur le bouton à droite *New* puis sélectionnez *Python 3*. Vous noterez au passage qu'il est également possible de créer un fichier texte, un répertoire ou bien encore de lancer un terminal (voir figure 18.3).



FIGURE 18.3 – Création d'un nouveau notebook

Le notebook créé ne contient qu'une cellule vide.

La première chose à faire est de donner un nom à votre notebook en cliquant sur *Untitled*, à droite du logo de Jupyter. Si le nom de votre notebook est *test* alors le fichier *test.ipynb* sera créé. L'extension *.ipynb* est l'extension de fichier des notebooks Jupyter.

Vous pouvez entrer une instruction Python dans la première cellule. Par exemple :

```
1 a = 2
2 b = 3
3 print(a+b)
```

Pour exécuter le contenu de cette cellule, vous avez plusieurs possibilités :

- Cliquer sur le menu *Cell*, puis *Run Cells*.

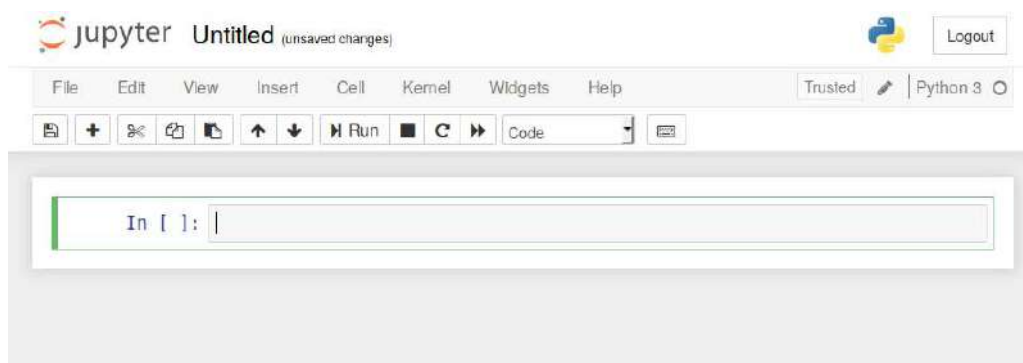


FIGURE 18.4 – Nouveau notebook

- Cliquer sur le bouton *Run* (sous la barre de menu).
- Appuyer simultanément sur les touches *Ctrl + Entrée*.

Dans tous les cas, vous devriez obtenir quelque chose qui ressemble à l'image 18.5. La notation `In [1]` à gauche de la cellule indique qu'il s'agit de la première cellule exécutée.

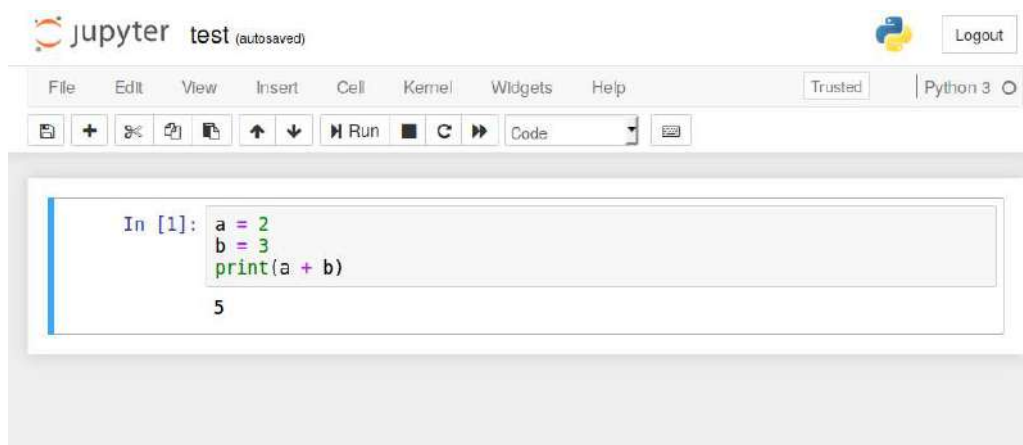


FIGURE 18.5 – Exécution d'une première cellule

Pour créer une nouvelle cellule, vous avez plusieurs possibilités ici encore :

- Cliquer sur l'icône *+* sous la barre de menu.
- Cliquer sur le menu *Insert*, puis *Insert Cell Below*.

Une nouvelle cellule vide devrait apparaître.

Vous pouvez également créer une nouvelle cellule en positionnant votre curseur dans la première cellule, puis en appuyant simultanément sur les touches *Maj + Entrée* (ou *Alt + Entrée*). Si vous utilisez cette combinaison de touches, vous remarquerez que le numéro à gauche de la première cellule est passée de `In [1]` à `In [2]` car vous avez exécuté la première cellule *puis* créez une nouvelle cellule.

Vous pouvez ainsi créer plusieurs cellules les unes à la suite des autres. Un objet créé dans une cellule antérieure sera disponible dans les cellules suivantes. Par exemple dans le figure 18.6, nous avons 4 cellules. Vous remarquerez que pour les cellules 3 et 4, le résultat renvoyé par le code Python est précédé par `Out [3]` et `Out [4]`.

Dans un notebook Jupyter, il est parfaitement possible de réexécuter une cellule précédente. Par exemple la première cellule qui porte désormais à sa gauche la numérotation `In [5]` (voir figure 18.7). Nous attirons votre attention que cette manipulation peut prêter à confusion. Si vous avez besoin de relancer une cellule antérieure, nous vous recommandons de relancer complètement l'exécution de

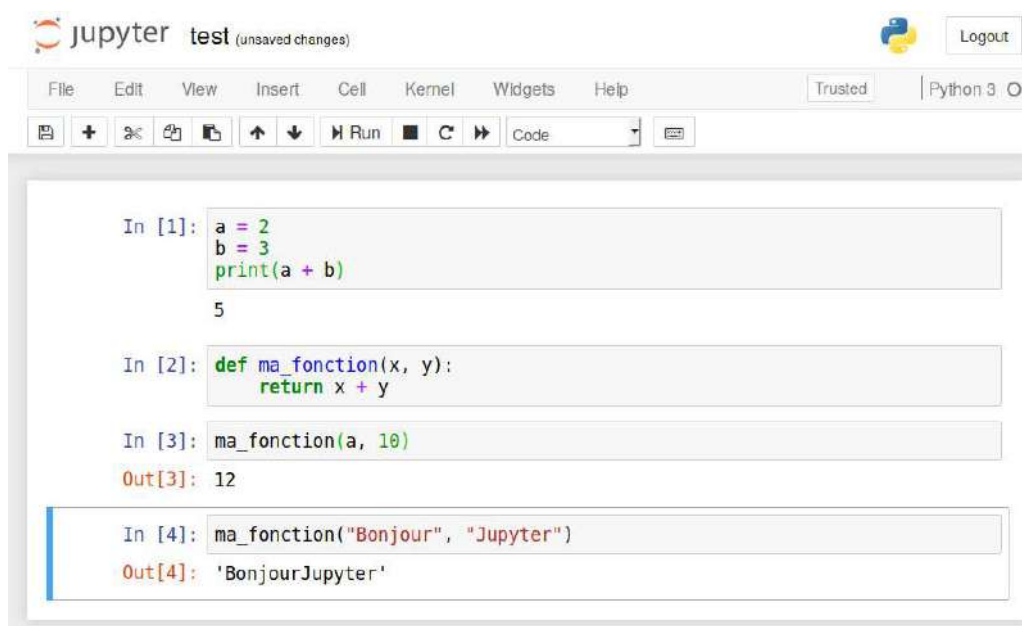


FIGURE 18.6 – Notebook avec plusieurs cellules de code Python

toutes les cellules, de la première à la dernière, en cliquant sur le menu *Kernel* puis *Restart & Run All* et enfin de valider le message *Restart and Run All Cells*.

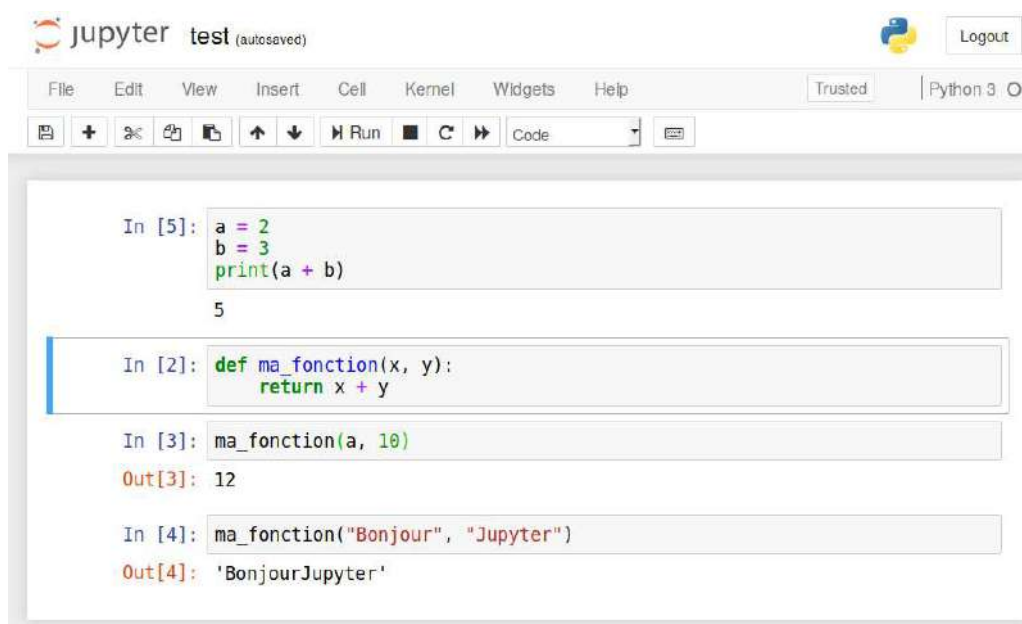


FIGURE 18.7 – Notebook avec une cellule ré-exécutée

18.2 Le format Markdown

Dans le tout premier exemple (figure 18.1), nous avons vu qu'il était possible de mettre du texte au format Markdown dans une cellule.

Il faut cependant indiquer à Jupyter que cette cellule est au format Markdown en cliquant sur *Code* sous la barre de menu puis en choisissant *Markdown*.

Le format Markdown permet de rédiger du texte formaté (gras, italique, lien, titres...) avec quelques

balises très simples. Voici un exemple dans une notebook Jupyter figure 18.8 et le rendu lorsque la cellule est exécutée (figure 18.9)

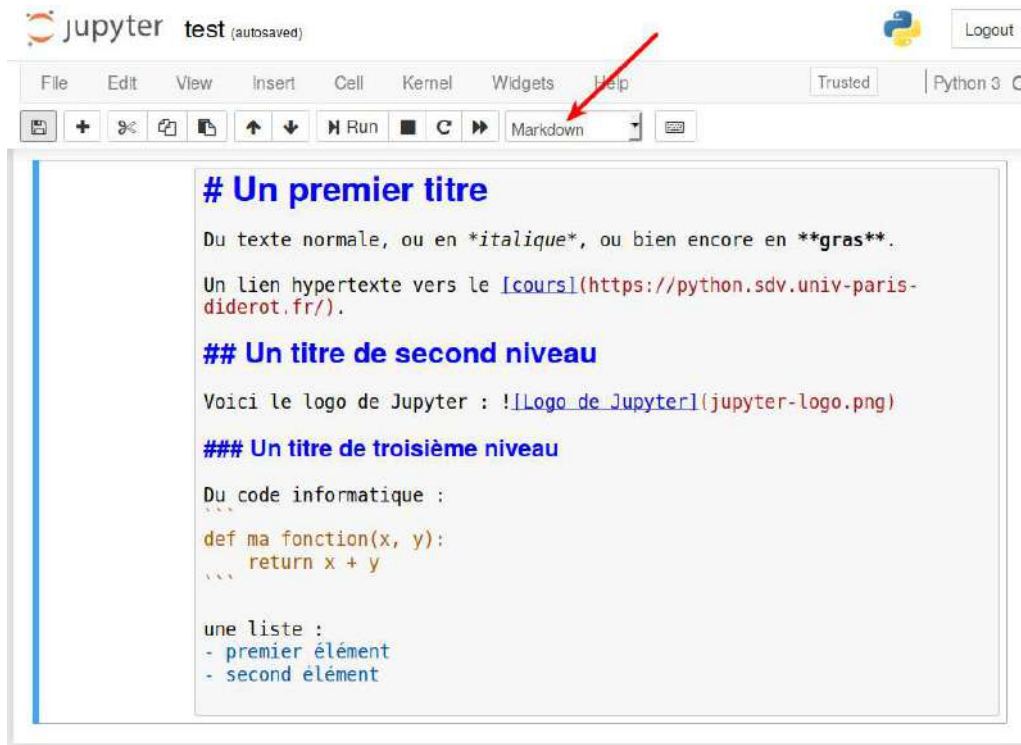


FIGURE 18.8 – Notebook avec une cellule au format Markdown

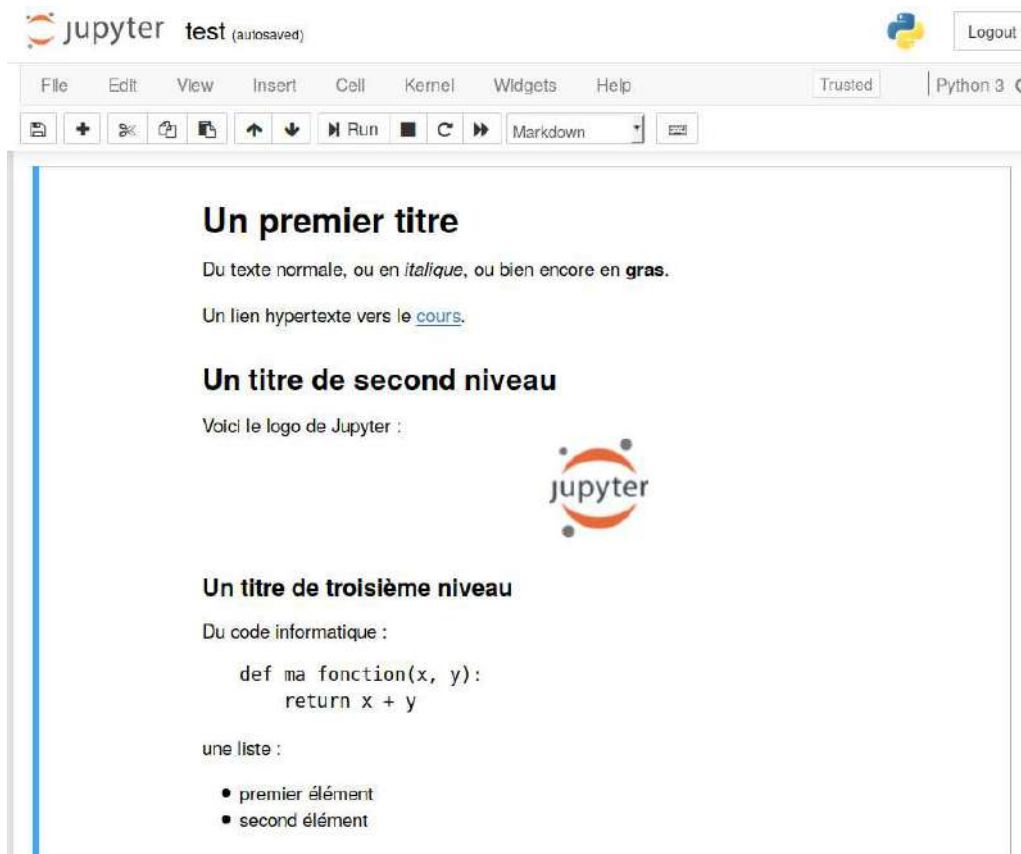


FIGURE 18.9 – Notebook avec une cellule au format Markdown (après exécution)

Notez qu'une cellule Markdown n'a pas les marqueurs In [] à sa gauche.

Le format Markdown permet de rapidement et très simplement rédiger du texte structuré. Ce cours est par exemple rédigé en Markdown ;-)

Nous vous conseillons d'explorer les possibilités du Markdown en consultant la page Wikipédia¹ ou directement la page de référence².

18.3 Des graphiques dans les notebooks

Un autre intérêt des notebooks Jupyter est de pouvoir y incorporer des graphiques réalisés avec la bibliothèque `matplotlib`.

Voici un exemple en reprenant un graphique présenté dans le chapitre sur les modules d'intérêt en bioinformatique (figure 18.10) :

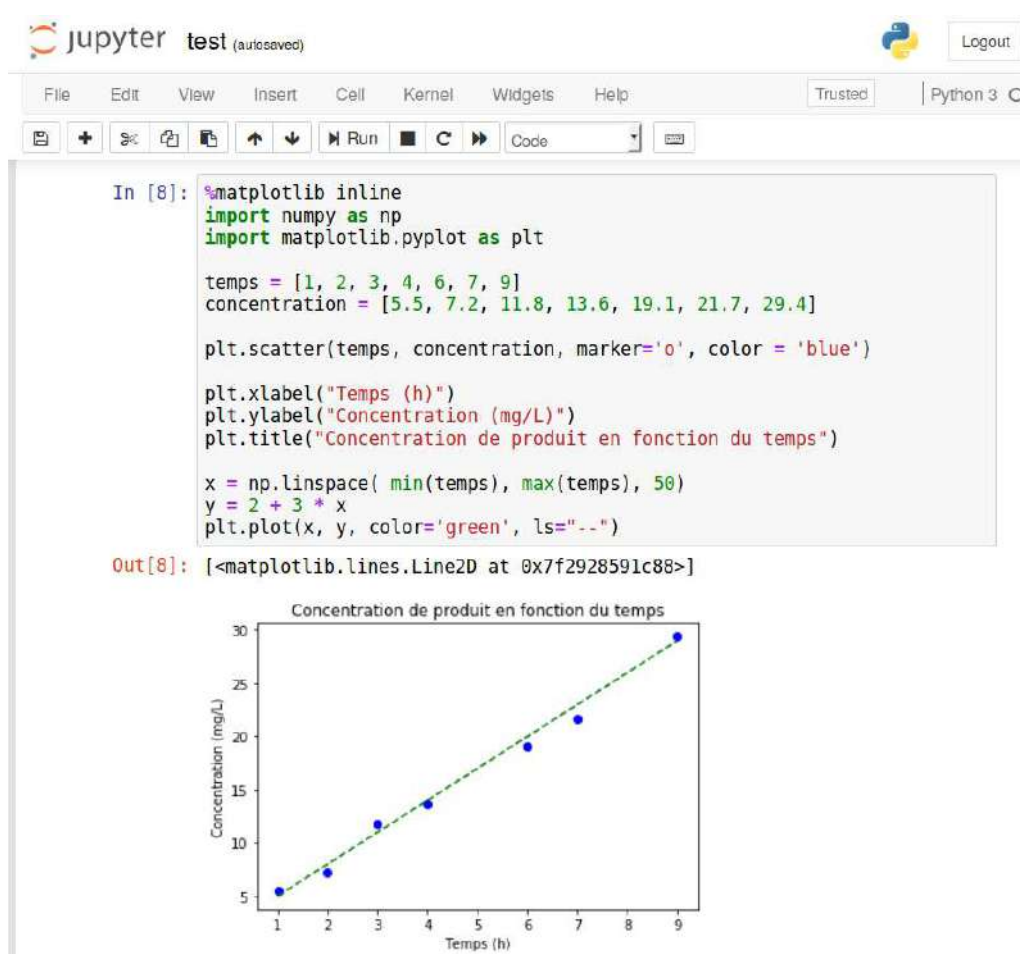


FIGURE 18.10 – Incorporation d'un graphique dans un notebook Jupyter

La différence notable est l'utilisation de la commande

```
1 %matplotlib inline
```

qui n'est à lancer qu'une seule fois (en générale dans la première cellule du notebook).

1. <https://fr.wikipedia.org/wiki/Markdown>

2. <https://daringfireball.net/projects/markdown/syntax>

18.4 Les *magic commands*

La commande précédente (`%matplotlib inline`) est une *magic command*. Il en existe beaucoup, en voici quelques unes :

`%whos` liste tous les objets (variables, fonctions, modules...) utilisés dans le notebook. Par exemple :

```
In [9]: %whos
```

Variable	Type	Data/Info
a	int	2
b	int	3
concentration	list	n=7
ma_fonction	function	<function ma_fonction at 0x7f2950b46048>
np	module	<module 'numpy' from '/ho<...>kages/numpy/_init__.py'>
plt	module	<module 'matplotlib.pyplot' from '/es/matplotlib/pyplot.py'>
temps	list	n=7
x	ndarray	50: 50 elems, type 'float64', 400 bytes
y	ndarray	50: 50 elems, type 'float64', 400 bytes

FIGURE 18.11 – *Magic command %whos*

`%history` liste toutes les commandes Python lancées dans un notebook. Par exemple :

```
In [7]: %history
```

```
a = 2
b = 3
print(a + b)
def ma_fonction(x, y):
    return x + y
ma_fonction(a, 10)
ma_fonction("Bonjour", "Jupyter")
%matplotlib inline
import numpy as np
import matplotlib.pyplot as plt

temps = [1, 2, 3, 4, 6, 7, 9]
concentration = [5.5, 7.2, 11.8, 13.6, 19.1, 21.7, 29.4]

plt.scatter(temps, concentration, marker='o', color = 'blue')

plt.xlabel("Temps (h)")
plt.ylabel("Concentration (mg/L)")
plt.title("Concentration de produit en fonction du temps")

x = np.linspace( min(temps), max(temps), 50)
y = 2 + 3 * x
plt.plot(x, y, color='green', ls="--")
%whos
%history
```

FIGURE 18.12 – *Magic command %history*

Enfin, il est possible de lancer une commande Unix depuis un notebook Jupyter. il faut pour cela précéder la commande du symbole `!`. Exemple avec la commande `ls` qui affiche le contenu d'un répertoire :

```
In [12]: !ls
```

```
jupyter-exemple.ipynb  jupyter-logo.png  test.ipynb
```

FIGURE 18.13 – Commande `!ls`

Chapitre 19

Avoir la classe avec les objets

Une classe permet de définir des objets qui sont des représentants (des **instances**) de cette classe. Les objets peuvent posséder des **attributs** (variables associées aux objets) et des **méthodes** (qui peuvent être vues comme des fonctions associées aux objets).

19.1 Construction d'une classe

Exemple de la classe Rectangle :

```
1 class Rectangle:
2     """
3     Ceci est la classe Rectangle
4     """
5
6     def __init__(self, long=0.0, larg=0.0, coul="blanc"):
7         """
8         Initialisation d'un objet.
9
10        Définition des attributs avec des valeurs par défaut.
11        """
12        self.longueur = long
13        self.largeur = larg
14        self.couleur = coul
15
16
17    def calcule_surface(self):
18        """
19        Méthode qui calcule la surface.
20        """
21        return self.longueur * self.largeur
22
23
24    def change_carre(self, cote):
25        """
26        Méthode qui transforme un rectangle en carré.
27        """
28        self.longueur = cote
29        self.largeur = cote
```

Dans cet exemple, `longueur`, `largeur` et `couleur` sont des attributs alors que `calcule_surface()` et `change_carre()` sont des méthodes. Tous les attributs et toutes les méthodes se réfèrent toujours à `self` qui désigne l'objet lui-même. Attention, les méthodes prennent au moins `self` comme argument.

19.2 Utilisation de la classe Rectangle

Création d'un objet `Rectangle` avec les paramètres par défaut :

```
1 rect1 = Rectangle()
2 print("longueur {} / largeur {} / couleur {}".format( rect1.longueur, rect1.largeur, rect1.couleur ))
3 print("surface = {:.2f} m2".format( rect1.calcule_surface() ))
```

Résultat :

```
1 longueur 0.0 / largeur 0.0 / couleur blanc
2 surface = 0.00 m2
```

19.3 Modification du rectangle en carré

```
1 rect1.change_carre(30)
2 rect1.calcule_surface()
```

Résultat :

```
1 surface = 900.00 m2
```

19.4 Création d'un objet Rectangle avec des paramètres imposés

```
1 rect2 = Rectangle(2, 3, "rouge")
2 print("longueur {} / largeur {} / couleur {}".format( rect2.longueur, rect2.largeur, rect2.couleur ))
3 print("surface = {:.2f} m2".format( rect2.calcule_surface() ))
```

Résultat :

```
1 longueur 2 / largeur 3 / couleur rouge
2 surface = 6.00 m2
```

19.5 Exercices

Conseil : pour ces exercices, écrivez des scripts dans des fichiers, puis exécutez-les dans un *shell*.

19.5.1 Classe Rectangle

Entraînez-vous avec la classe `Rectangle`. Créez la méthode `calcule_perimetre()` qui calcule le périmètre d'un objet rectangle.

19.5.2 Classe Atome

Créez une nouvelle classe `Atome` avec les attributs `x`, `y`, `z` (qui contiennent les coordonnées atomiques) et la méthode `calcul_distance()` qui calcule la distance entre deux atomes. Testez cette classe sur plusieurs exemples.

19.5.3 Classe Atome améliorée

Améliorez la classe `Atome` en lui ajoutant un nouvel attribut `masse` qui correspond à la masse atomique et `numero_atomique`) et une nouvelles méthodes `calculer_centre_masse`.

Chapitre 20

Pour aller plus loin

20.1 Shebang et `/usr/bin/env python3`

Lorsque vous programmez sur un système Unix, le shebang¹ correspond aux caractères `#!` qui se trouve au début de la première ligne d'un script. Le shebang est suivi du chemin complet du programme qui interprète le script.

En Python, on trouve souvent la notation

```
1 #! /usr/bin/python3
```

Cependant, l'exécutable `python3` ne se trouve pas toujours dans le répertoire `/usr/bin`. Pour maximiser la portabilité de votre script Python sur plusieurs systèmes Unix, utilisez plutôt cette notation :

```
1 #! /usr/bin/env python3
```

Dans le cas présent, on appelle le programme d'environnement `env` (qui se situe toujours dans le répertoire `/usr/bin`) pour lui demander où se trouve l'exécutable `python3`.

20.2 Différences Python 2 et Python 3

20.2.1 Interpréteur par défaut

Par défaut sous Ubuntu 16.04, l'exécutable `python` va lancer Python 2 :

```
1 $ python
2 Python 2.7.12 (default, Nov 19 2016, 06:48:10)
3 [GCC 5.4.0 20160609] on linux2
4 Type "help", "copyright", "credits" or "license" for more information.
5 >>>
```

Pensez à explicitement préciser la version de Python qui vous intéresse, ici Python 3 :

```
1 $ python3
2 Python 3.5.2 (default, Nov 17 2016, 17:05:23)
3 [GCC 5.4.0 20160609] on linux
4 Type "help", "copyright", "credits" or "license" for more information.
5 >>>
```

1. <http://fr.wikipedia.org/wiki/Shebang>

20.2.2 La fonction print()

La fonction `print()` en Python 2 s'utilise sans parenthèse. Par exemple :

```
1 >>> print 12
2 12
3 >>> print "girafe"
4 girafe
```

Par contre en Python 3, si vous n'utilisez pas de parenthèse, Python vous renverra une erreur :

```
1 >>> print 12
2 File "<stdin>", line 1
3   print 12
4     ^
5 SyntaxError: Missing parentheses in call to 'print'
```

20.2.3 Division d'entiers

En Python 3, la division de deux entiers, se fait *naturellement*, c'est-à-dire que l'opérateur `/` renvoie systématiquement un *float*. Par exemple :

```
1 >>> 3 / 4
2 0.75
```

Il est également possible de réaliser une division entière avec l'opérateur `//` :

```
1 >>> 3 // 4
2 0
```

La division entière renvoie finalement la partie entière du nombre 0.75, c'est à dire 0.

Attention! En Python 2, la division de deux entiers avec l'opérateur `/` correspond à la division entière, c'est-à-dire le résultat arrondi à l'entier inférieur. Par exemple :

```
1 >>> 3 / 5
2 0
3 >>> 4 / 3
4 1
```

Faites très attention à cet aspect si vous programmez encore en Python 2, c'est une source d'erreur récurrente.

20.2.4 La fonction range()

En Python 3, la fonction `range()` est un générateur, c'est-à-dire que cette fonction va itérer sur le nombre entier donné en argument. On ne peut pas l'utiliser seule :

```
1 >>> range(3)
2 range(0, 3)
```

Lorsqu'on l'utilise dans une boucle `for`, `range(3)` va produire successivement les nombres 0, 1 puis 2. Par exemple :

```

1 >>> for i in range(3):
2     ...     print(i)
3     ...
4     0
5     1
6     2

```

En Python 2, la fonction `range()` renvoie une liste. Par exemple :

```

1 >>> range(3)
2 [0, 1, 2]
3 >>> range(2, 6)
4 [2, 3, 4, 5]

```

La création de liste avec `range()` était pratique mais très peu efficace en mémoire lorsque l'argument de `range()` était un grand nombre.

D'ailleurs la fonction `xrange()` est disponible en Python 2 pour faire la même chose que la fonction `range()` en Python 3. Attention, ne vous mélangez pas les pinceaux !

```

1 >>> range(3)
2 [0, 1, 2]
3 >>> xrange(3)
4 xrange(3)

```

Remarque : pour générer une liste d'entiers avec la fonction `range()` en Python 3, vous avez vu dans le chapitre sur les *Listes* qu'il suffit de l'associer avec la fonction `list()`. Par exemple :

```

1 >>> list(range(10))
2 [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]

```

20.2.5 Encodage et utf-8

En Python 3, vous pouvez utiliser des caractères accentués dans les commentaires ou dans les chaînes de caractères.

Ce n'est pas le cas en Python 2. Si un caractère accentué est présent dans votre code, cela occasionnera une erreur de ce type lors de l'exécution de votre script :

```

1 SyntaxError: Non-ASCII character '\xc2' in file xxx on line yyy, but no encoding declared; see

```

Pour éviter ce genre de désagrément, ajoutez la ligne suivante à la deuxième ligne (la position est importante²) de votre script :

```

1 # coding: utf-8

```

En résumé, tous vos scripts en Python 2 devraient ainsi débiter par les lignes :

```

1 #! /usr/bin/env python2
2 # coding: utf-8

```

Remarque : l'encodage utf-8 peut aussi être déclaré de cette manière :

2. <http://www.python.org/dev/peps/pep-0263/>

```
1 # -*- coding: utf-8 -*-
```

mais c'est un peu plus long à écrire.

20.3 Liste de compréhension

Une manière originale et très puissante de générer des listes est la compréhension de listes. Pour plus de détails, consultez à ce sujet le site de Python³ et celui de Wikipédia⁴.

Voici quelques exemples.

20.3.1 Nombres pairs compris entre 0 et 30

```
1 >>> print([i for i in range() if i%2 == 0])
2 [0, 2, 4, 6, 8, 10, 12, 14, 16, 18, 20, 22, 24, 26, 28, 30]
```

20.3.2 Jeu sur la casse des mots d'une phrase

```
1 >>> message = "C'est sympa la BioInfo"
2 >>> msg_lst = message.split()
3 >>> print([[m.upper(), len(m)] for m in msg_lst])
4 [['C'EST', 5], ['SYMPA', 5], ['LA', 2], ['BIOINFO', 7]]
```

20.3.3 Formatage d'une séquence avec 60 caractères par ligne

```
1 # exemple d'une séquence de 150 alanines
2 >>> seq = "A"*150
3 >>> width= 60
4 >>> print("\n".join( [seq[i:i+width] for i in range(0,len(seq),width)] ))
5 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
6 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
7 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
```

20.3.4 Formatage fasta d'une séquence (avec la ligne de commentaire)

Exemple d'une séquence constituée de 150 alanines :

```
1 >>> com = "séquence de 150 alanines"
2 >>> seq = "A"*150
3 >>> width = 60
4 >>> print(">"+com+"\n"+"".join( [seq[i:i+width] for i in range(0,len(seq),width)] ))
5 >séquence de 150 alanines
6 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
7 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
8 AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
```

3. <http://www.python.org/dev/peps/pep-0202/>

4. http://fr.wikipedia.org/wiki/Comprehension_de_liste

20.3.5 Sélection des carbones alpha dans un fichier pdb

Exemple avec la structure de la barstar⁵ :

```

1 >>> with open("1BTA.pdb", "r") as fichier:
2     ...     CA_lines = [l for l in fichier if l.split()[0] == "ATOM" and l.split()[2] == "CA"]
3     ...
4 >>> print(len(CA_lines))
5 89

```

20.4 Gestion des erreurs

La gestion des erreurs permet d'éviter que votre programme plante en prévoyant vous même les sources d'erreurs éventuelles.

Voici un exemple dans lequel on demande à l'utilisateur d'entrer un nombre entier, puis on affiche ce nombre.

```

1 >>> nb = int(input("Entrez un nombre entier : "))
2 Entrez un nombre entier : 23
3 >>> print(nb)
4 23

```

La fonction `input()` demande à l'utilisateur de saisir une chaîne de caractères. Cette chaîne de caractères est ensuite transformée en nombre entier avec la fonction `int()`.

Si l'utilisateur ne rentre pas un nombre, voici ce qui se passe :

```

1 >>> nb = int(input("Entrez un nombre entier : "))
2 Entrez un nombre entier : ATCG
3 Traceback (most recent call last):
4   File "<stdin>", line 1, in <module>
5 ValueError: invalid literal for int() with base 10: 'ATCG'

```

L'erreur provient de la fonction `int()` qui n'a pas pu convertir la chaîne de caractères "ATCG" en nombre entier, ce qui est parfaitement normal.

Le jeu d'instructions `try / except` permet de tester l'exécution d'une commande et d'intervenir en cas d'erreur.

```

1 >>> try:
2     ...     nb = int(input("Entrez un nombre entier : "))
3     ... except:
4     ...     print("Vous n'avez pas entré un nombre entier !")
5     ...
6 Entrez un nombre entier : ATCG
7 Vous n'avez pas entré un nombre entier !

```

Dans cette exemple, l'erreur renvoyée par la fonction `int()` (qui ne peut pas convertir "ATCG" en nombre entier) est interceptée et déclenche l'affichage du message d'avertissement.

On peut ainsi redemander sans cesse un nombre entier à l'utilisateur, jusqu'à ce que celui-ci en rentre bien un.

5. <http://www.rcsb.org/pdb/explore.do?structureId=1BTA>

```

1 >>> while True:
2     ...     try:
3     ...         nb = int(input("Entrez un nombre entier : "))
4     ...         print("Le nombre est", nb)
5     ...         break
6     ...     except:
7     ...         print("Vous n'avez pas entré un nombre entier !")
8     ...         print("Essayez encore")
9     ...
10 Entrez un nombre entier : ATCG
11 Vous n'avez pas entré un nombre entier !
12 Essayez encore
13 Entrez un nombre entier : toto
14 Vous n'avez pas entré un nombre entier !
15 Essayez encore
16 Entrez un nombre entier : 3.2
17 Vous n'avez pas entré un nombre entier !
18 Essayez encore
19 Entrez un nombre : 55
20 Le nombre est 55

```

Notez que dans cet exemple, l'instruction `while 1` est une boucle infinie (car la condition `True` est toujours vérifiée) dont l'arrêt est forcé par la commande `break` lorsque l'utilisateur a effectivement bien rentré un nombre entier.

La gestion des erreurs est très utile dès lors que des données extérieures entrent dans un programme Python, que ce soit directement par l'utilisateur (avec la fonction `input()`) ou par des fichiers.

Vous pouvez par exemple vérifier qu'un fichier a bien été ouvert.

```

1 >>> nom = "toto.pdb"
2 >>> try:
3     ...     with open(nom, "r") as fichier:
4     ...         for ligne in fichier:
5     ...             print(ligne)
6     ... except:
7     ...     print("Impossible d'ouvrir le fichier", nom)

```

Si une erreur est déclenchée, c'est sans doute que le fichier n'existe pas à l'emplacement indiqué sur le disque ou que vous n'avez pas les droits pour le lire.

Il est également possible de spécifier le type d'erreur à gérer. Le premier exemple que nous avons étudié peut s'écrire :

```

1 >>> try:
2     ...     nb = int(input("Entrez un nombre entier : "))
3     ... except ValueError:
4     ...     print("Vous n'avez pas entré un nombre entier !")
5     ...
6 Entrez un nombre entier : ATCG
7 Vous n'avez pas entré un nombre entier !

```

Ici, on intercepte une erreur de type `ValueError`, ce qui correspond bien à un problème de conversion avec `int()`. Il existe d'autres types d'erreurs comme `RuntimeError`, `TypeError`, `NameError`, `IOError`,

etc.

Enfin, on peut aussi être très précis dans le message d'erreur. Observez la fonction `download_page()` qui, avec le module `urllib`, télécharge un fichier sur internet.

```

1 import urllib.request
2
3 def download_page(address):
4     error = ""
5     page = ""
6     try:
7         data = urllib.request.urlopen(address)
8         page = data.read()
9     except IOError as e:
10        if hasattr(e, 'reason'):
11            error = "Cannot reach web server: " + str(e.reason)
12        if hasattr(e, 'code'):
13            error = "Server failed {:d}".format(e.code)
14    return page, error
15
16 data, error = download_page("https://files.rcsb.org/download/1BTA.pdb")
17
18 if error:
19     print("Erreur rencontrée : {}".format(error))
20 else:
21     with open("proteine.pdb", "w") as prot:
22         prot.write(data.decode('utf-8'))
23     print("Protéine enregistrée")

```

La variable `e` est une instance (un représentant) de l'erreur de type `IOError`. Certains de ces attributs sont testés avec la fonction `hasattr()` pour ainsi affiner le message renvoyé (ici contenu dans la variable `error`).

Si tout se passe bien, la page est téléchargée et stockée dans la variable `data`, puis ensuite enregistrée sur le disque.

20.5 Sauvegardez votre historique de commandes

Vous pouvez sauvegarder l'historique des commandes utilisées dans l'interpréteur Python avec le module `readline`.

```

1 >>> print("hello")
2 hello
3 >>> a = 22
4 >>> a = a + 11
5 >>> print(a)
6 33
7 >>> import readline
8 >>> readline.write_history_file()

```

Quittez Python. L'historique de toutes vos commandes est dans votre répertoire personnel, dans le fichier `.history`.

Relancez l'interpréteur Python.

```
1 >>> import readline
2 >>> readline.read_history_file()
```

Vous pouvez accéder aux commandes de la session précédente avec la flèche du haut de votre clavier. D'abord les commandes `readline.read_history_file()` et `import readline` de la session actuelle, puis `print(a)`, `a = a + 11`, `a = 22...`

Chapitre 21

Mini-projets

Dans ce chapitre, nous vous proposons quelques scénarios pour développer vos compétences en Python et mettre en oeuvre les concepts que vous avez rencontrés dans les chapitres précédents.

21.1 Mots anglais dans le protéome humain

L'objectif de ce premier projet est de découvrir si des mots anglais peuvent se retrouver dans les séquences du protéome humain, c'est-à-dire dans les séquences de l'ensemble des protéines humaines.

21.1.1 Composition aminée

Dans un premier temps, composez 5 mots anglais avec les 20 acides aminés.

21.1.2 Des mots

Téléchargez le fichier `english-common-words.txt`¹. Ce fichier contient les 3000 mots anglais les plus fréquents, à raison d'1 mot par ligne.

Créez un script `words-in-proteome.py` et écrivez la fonction `read_words()` qui va lire les mots contenus dans le fichier dont le nom est fourni en argument du script et renvoyer une liste contenant les mots convertis en majuscule et composés de 3 caractères ou plus.

Dans le programme principal, affichez le nombre de mots sélectionnés.

21.1.3 Des protéines

Téléchargez maintenant le fichier `human-proteome.fasta`². Attention, ce fichier est assez gros. Ce fichier provient de la banque de données UniProt à partir de cette page³.

Voici les premières lignes de ce fichier ([...] indique une coupure que nous avons faite) :

```
1 >sp|095139|NDUB6_HUMAN NADH dehydrogenase [ubiquinone] 1 beta [...]
2 MTGYTPDEKLRLLQLRELRRRWLKDQELSPREPVLPQKMGPMKFWNKFLENKSPWRKM
3 VHGVIYKKSIFVFTHVLVPVWIIHYMKYHVSEKPYGIVEKKSRIFFPGDTILETGEVIPP
4 KEFPDQHH
```

1. <https://python.sdv.univ-paris-diderot.fr/data-files/english-common-words.txt>
2. <https://python.sdv.univ-paris-diderot.fr/data-files/human-proteome.fasta>
3. https://www.uniprot.org/help/human_proteome

```

5 >sp|075438|NDUB1_HUMAN NADH dehydrogenase [ubiquinone] 1 beta [...]
6 MVNLLQIVRDHWHVHVLVPMGFVIGCYLDRKSDERLTAFRNKSMFLKRELQPSEEVTK
7 >sp|Q8N4C6|NIN_HUMAN Ninein OS=Homo sapiens OX=9606 GN=NIN PE=1 SV=4
8 MDEVEQDQHEARLKELFDSFDTTGTGSLGQEELTDLCHMLSLEEVAPVLQQTLLQDNLLG
9 RVHFDQFKEALILILSRTLSEEHFQEPDCSLEAQPKYVRGGKRYGRRSLPEFQESVEEF
10 PEVTVIEPLDEEARPSHIPAGDCSEHWKTQRSEEYEAEGQLRFWNPDDLNASQSGSSPPQ

```

Toujours dans le script `words-in-proteome.py`, écrivez la fonction `read_sequences()` qui va lire le protéome dans le fichier dont le nom est fourni en second argument du script. Cette fonction va renvoyer un dictionnaire dont les clefs sont les identifiants des protéines (par exemple, `095139`, `075438`, `Q8N4C6`) et dont les valeurs associées sont les séquences.

Dans le programme principal, affichez le nombre de séquences lues. À des fins de test, affichez également la séquence associée à la protéine `095139`.

21.1.4 À la pêche aux mots

Écrivez maintenant la fonction `search_words_in_proteome()` qui prend en argument la liste de mots et le dictionnaire contenant les séquences des protéines et qui va compter le nombre de séquences dans lesquelles un mot est présent. Cette fonction renverra un dictionnaire dont les clefs sont les mots et les valeurs le nombre de séquences qui contiennent ces mots. La fonction affichera également le message suivant pour les mots trouvés dans le protéome :

```

1 ACCESS found in 1 sequences
2 ACID found in 38 sequences
3 ACT found in 805 sequences
4 [...]

```

Cette étape prend quelques minutes. Soyez patient.

21.1.5 Et le mot le plus fréquent est...

Pour terminer, écrivez maintenant la fonction `find_most_frequent_word()` qui prend en argument le dictionnaire renvoyé par la précédente fonction `search_words_in_proteome()` et qui affiche le mot trouvé dans le plus de protéines, ainsi que le nombre de séquences dans lesquelles il a été trouvé, sous la forme :

```

1 => xxx found in yyy sequences

```

Quel est ce mot ?

Quel pourcentage des séquences du protéome contiennent ce mot ?

21.1.6 Pour être plus complet

Jusqu'à présent, nous avons déterminé, pour chaque mot, le nombre de séquences dans lesquelles il apparaissait. Nous pourrions aller plus loin et calculer aussi le nombre de fois que chaque mot apparaît dans les séquences.

Pour cela modifier la fonction `search_words_in_proteome()` de façon à compter le nombre d'occurrences d'un mot dans les séquences. La méthode `.count()` vous sera utile.

Déterminez alors quel mot est le plus fréquent dans le protéome humain.

21.2 genbank2fasta

Ce projet consiste à écrire un convertisseur de fichier, du format GenBank au format fasta.

Pour cela, nous allons utiliser le fichier genbank du chromosome I de la levure du boulanger *Saccharomyces cerevisiae*. Vous pouvez télécharger ce fichier :

- soit via le lien sur le site du cours NC_001133.gbk⁴ ;
- soit directement sur la page de *Saccharomyces cerevisiae* S288c chromosome I, complete sequence⁵ sur le site du NCBI, puis en cliquant sur *Send to*, puis *Complete Record*, puis *Choose Destination : File*, puis *Format : GenBank (full)* et enfin sur le bouton *Create File*.

Vous pouvez consulter les caractéristiques des formats fasta et genbank dans l'annexe *Quelques formats de données rencontrés en biologie*.

Dans la suite, nous vous proposons deux manières de procéder, avec et sans expression régulière selon si vous avez ou non lu et acquis les expressions régulières (Chapitre 15).

21.2.1 genbank2fasta sans expression régulière

Si vous n'avez pas encore travaillé les expressions régulières (Chapitre 15), vous êtes au bon endroit ! Ces fameuses expressions régulières permettent de traiter ce problème de manière puissante et élégante, mais il est tout à fait possible de réaliser ce mini projet sans elles.

Lecture du fichier

Créez un script `genbank2fasta.py` et créez la fonction `lit_fichier()` qui prend en argument le nom du fichier et qui renvoie le contenu du fichier sous forme d'une liste de lignes, chaque ligne étant elle-même une chaîne de caractères.

Testez cette fonction avec le fichier GenBank NC_001133.gbk et affichez le nombre de lignes lues.

Extraction du nom de l'organisme

Dans le même script, ajoutez la fonction `extrait_organisme()` qui prend en argument le contenu du fichier précédemment obtenu avec la fonction `lit_fichier()` (sous la forme d'une liste de lignes) et qui renvoie le nom de l'organisme. Pour récupérer la bonne ligne vous pourrez tester si les premiers caractères de la ligne contiennent le mot-clé `ORGANISM`.

Testez cette fonction avec le fichier GenBank NC_001133.gbk et affichez le nom de l'organisme.

Recherche des gènes

Dans le fichier GenBank, les gènes sens sont notés de cette manière :

1

gene	58..272
------	---------

ou

1

gene	<2480..>2707
------	--------------

4. https://python.sdv.univ-paris-diderot.fr/data-files/NC_001133.gbk

5. https://www.ncbi.nlm.nih.gov/nuccore/NC_001133

et les gènes antisens (ou encore complémentaires) de cette façon :

```
1 gene complement(55979..56935)
```

ou

```
1 gene complement(<13363..>13743)
```

Les valeurs numériques séparées par .. indiquent la position du gène dans le génome (numéro de la première base, numéro de la dernière base).

Remarque : le symbole < indique un gène partiel sur l'extrémité 5', c'est-à-dire que le codon START correspondant est incomplet. Respectivement, le symbole > désigne un gène partiel sur l'extrémité 3', c'est-à-dire que le codon STOP correspondant est incomplet. Pour plus de détails, consultez la documentation du NCBI sur les délimitations des gènes⁶. Nous vous proposons ici d'ignorer ces symboles > et <.

Repérez ces différents gènes dans le fichier `NC_001133.gbk`. Pour récupérer ces lignes de gènes il suffira donc de tester si la ligne commence par

```
1 gene
```

(c'est à dire 5 espaces, suivi du mot `gene`, suivi de 12 espaces). Pour savoir s'il s'agit d'un gène sur le brin direct ou complémentaire, il suffira de tester la présence du mot `complement` dans la ligne lue.

Ensuite si vous souhaitez récupérer la position de début et de fin de gène, nous vous conseillons d'utiliser la fonction `replace()` et de ne garder que les chiffres et les . Par exemple

```
1 gene <2480..>2707
```

sera transformé en

```
1 2480..2707
```

Enfin, avec la fonction `split()` vous pourrez facilement récupérer les deux entiers de début et de fin de gène.

Dans le même script `genbank2fasta.py`, ajoutez la fonction `recherche_genes()` qui prend en argument le contenu du fichier (sous la forme d'une liste de lignes) et qui renvoie la liste des gènes.

Chaque gène sera lui-même une liste contenant le numéro de la première base, le numéro de la dernière base et une chaîne de caractère "`sens`" pour un gène sens et "`antisens`" pour un gène antisens.

Testez cette fonction avec le fichier GenBank `NC_001133.gbk` et affichez le nombre de gènes trouvés, ainsi que le nombre de gènes sens et antisens.

Extraction de la séquence nucléique du génome

La taille du génome est indiqué sur la première ligne d'un fichier GenBank. Trouvez la taille du génome stocké dans le fichier `NC_001133.gbk`.

Dans un fichier GenBank, la séquence du génome se trouve entre les lignes

```
1 ORIGIN
```

et

6. <https://www.ncbi.nlm.nih.gov/Sitemap/samplerecord.html#BaseSpanB>


```
1 //
```

Trouvez dans le fichier `NC_001133.gbk` la première et dernière ligne de la séquence du génome.

Pour récupérer les lignes contenant la séquence, nous vous proposons d'utiliser un algorithme avec un drapeau `is_dnaseq` (qui vaudra `True` ou `False`). Voici l'algorithme proposé en pseudo-code :

```
1 is_dnaseq <- False
2 Lire chaque ligne du fichier gbk
3   si la ligne contient "\\\"
4     is_dnaseq <- False
5   si is_dnaseq vaut True
6     accumuler la séquence
7   si la ligne contient "ORIGIN"
8     is_dnaseq <- True
```

Au début ce drapeau aura la valeur `False`. Ensuite, quand il se mettra à `True`, on pourra lire les lignes contenant la séquence, puis quand il se remettra à `False` on arrêtera.

Une fois la séquence récupérée, il suffira d'éliminer les chiffres, retours chariots et autres espaces (*Conseil* : calculer la longueur de la séquence et comparer la à celle indiquée dans le fichier `gbk`).

Toujours dans le même script `genbank2fasta.py`, ajoutez la fonction `extrait_séquence()` qui prend en argument le contenu du fichier (sous la forme de liste de lignes) et qui renvoie la séquence nucléique du génome (dans une chaîne de caractères). La séquence ne devra pas contenir d'espaces, ni de chiffres ni de retours chariots.

Testez cette fonction avec le fichier GenBank `NC_001133.gbk` et affichez le nombre de bases de la séquence extraite. Vérifiez que vous n'avez pas fait d'erreur en comparant la taille de la séquence extraite avec celle que vous avez trouvée dans le fichier GenBank.

Construction d'une séquence complémentaire inverse

Toujours dans le même script, ajoutez la fonction `construit_comp_inverse()` qui prend en argument une séquence d'ADN sous forme de chaîne de caractères et qui renvoie la séquence complémentaire inverse (également sous la forme d'une chaîne de caractères).

On rappelle que construire la séquence complémentaire inverse d'une séquence d'ADN consiste à :

- Prendre la séquence complémentaire. C'est-à-dire à remplacer la base `a` par la base `t`, `t` par `a`, `c` par `g` et `g` par `c`.
- Prendre l'inverse. C'est-à-dire à que la première base de la séquence complémentaire devient la dernière base et réciproquement, la dernière base devient la première.

Pour vous faciliter le travail, ne travaillez que sur des séquences en minuscule.

Testez cette fonction avec les séquences `atcg`, `AATTCCGG` et `gattaca`.

Écriture d'un fichier fasta

Toujours dans le même script, ajoutez la fonction `ecrit_fasta()` qui prend en argument un nom de fichier (sous forme de chaîne de caractères), un commentaire (sous forme de chaîne de caractères) et une séquence (sous forme de chaîne de caractères) et qui écrit un fichier fasta. La séquence sera écrite sur des lignes ne dépassant pas 80 caractères.

Pour rappel, un fichier fasta suit le format suivant :

Lecture du fichier

Créez un script `genbank2fasta.py` et créez la fonction `lit_fichier()` qui prend en argument le nom du fichier et qui renvoie le contenu du fichier sous forme d'une liste de lignes, chaque ligne étant elle-même une chaîne de caractères.

Testez cette fonction avec le fichier `GenBank NC_001133.gbk` et affichez le nombre de lignes lues.

Extraction du nom de l'organisme

Dans le même script, ajoutez la fonction `extrait_organisme()` qui prend en argument le contenu du fichier précédemment obtenu avec la fonction `lit_fichier()` (sous la forme d'une liste de lignes) et qui renvoie le nom de l'organisme. Utilisez de préférence une expression régulière.

Testez cette fonction avec le fichier `GenBank NC_001133.gbk` et affichez le nom de l'organisme.

Recherche des gènes

Dans le fichier GenBank, les gènes sens sont notés de cette manière :

```
1 gene          58..272
```

ou

```
1 gene          <2480..>2707
```

et les gènes antisens de cette façon :

```
1 gene          complement(55979..56935)
```

ou

```
1 gene          complement(<13363..>13743)
```

Les valeurs numériques séparées par `..` indiquent la position du gène dans le génome (numéro de la première base, numéro de la dernière base).

Remarque : le symbole `<` indique un gène partiel sur l'extrémité 5', c'est-à-dire que le codon START correspondant est incomplet. Respectivement, le symbole `>` désigne un gène partiel sur l'extrémité 3', c'est-à-dire que le codon STOP correspondant est incomplet. Pour plus de détails, consultez la documentation du NCBI sur les délimitations des gènes⁷.

Repérez ces différents gènes dans le fichier `NC_001133.gbk`. Construisez deux expressions régulières pour extraire du fichier GenBank les gènes sens et les gènes antisens.

Modifiez ces expressions régulières pour que les numéros de la première et de la dernière base puissent être facilement extraits.

Dans le même script `genbank2fasta.py`, ajoutez la fonction `recherche_genes()` qui prend en argument le contenu du fichier (sous la forme d'une liste de lignes) et qui renvoie la liste des gènes.

Chaque gène sera lui-même une liste contenant le numéro de la première base, le numéro de la dernière base et une chaîne de caractère `"sens"` pour un gène sens et `"antisens"` pour un gène antisens.

Testez cette fonction avec le fichier `GenBank NC_001133.gbk` et affichez le nombre de gènes trouvés, ainsi que le nombre de gènes sens et antisens.

7. <https://www.ncbi.nlm.nih.gov/Sitemap/samplerecord.html#BaseSpanB>

Extraction de la séquence nucléique du génome

La taille du génome est indiqué sur la première ligne d'un fichier GenBank. Trouvez la taille du génome stocké dans le fichier `NC_001133.gbk`.

Dans un fichier GenBank, la séquence du génome se trouve entre les lignes

```
1 ORIGIN
```

et

```
1 //
```

Trouvez dans le fichier `NC_001133.gbk` la première et dernière ligne de la séquence du génome.

Construisez une expression régulière pour extraire du fichier GenBank les lignes correspondantes à la séquence du génome.

Modifiez ces expressions régulières pour que la séquence puisse être facilement extraite.

Toujours dans le même script, ajoutez la fonction `extrait_sequence()` qui prend en argument le contenu du fichier (sous la forme de liste de lignes) et qui renvoie la séquence nucléique du génome (dans une chaîne de caractères). La séquence ne devra pas contenir d'espaces.

Testez cette fonction avec le fichier GenBank `NC_001133.gbk` et affichez le nombre de bases de la séquence extraite. Vérifiez que vous n'avez pas fait d'erreur en comparant la taille de la séquence extraite avec celle que vous avez trouvée dans le fichier GenBank.

Construction d'une séquence complémentaire inverse

Toujours dans le même script, ajoutez la fonction `construit_comp_inverse()` qui prend en argument une séquence d'ADN sous forme de chaîne de caractères et qui renvoie la séquence complémentaire inverse (également sous la forme d'une chaîne de caractères).

On rappelle que construire la séquence complémentaire inverse d'une séquence d'ADN consiste à :

- Prendre la séquence complémentaire. C'est-à-dire à remplacer la base `a` par la base `t`, `t` par `a`, `c` par `g` et `g` par `c`.
- Prendre l'inverse. C'est-à-dire à que la première base de la séquence complémentaire devient la dernière base et réciproquement, la dernière base devient la première.

Pour vous faciliter le travail, ne travaillez que sur des séquences en minuscule.

Testez cette fonction avec les séquences `atcg`, `AATTCCGG` et `gattaca`.

Écriture d'un fichier fasta

Toujours dans le même script, ajoutez la fonction `ecrit_fasta()` qui prend en argument un nom de fichier (sous forme de chaîne de caractères), un commentaire (sous forme de chaîne de caractères) et une séquence (sous forme de chaîne de caractères) et qui écrit un fichier fasta. La séquence sera écrite sur des lignes ne dépassant pas 80 caractères.

Pour rappel, un fichier fasta suit le format suivant :

```
1 >commentaire
2 sequence sur une ligne de 80 caractères maxi
3 suite de la séquence .....
4 suite de la séquence .....
```


Chapitre 22

Annexe : Quelques formats de données rencontrés en biologie

22.1 Fasta

Le format Fasta est utilisé pour stocker un ou plusieurs séquences, d'ADN, d'ARN ou de protéines.

Ces séquences sont classiquement représentées sous la forme :

```
1 >en-tête
2 séquence avec un nombre maximum de caractères par ligne
3 séquence avec un nombre maximum de caractères par ligne
4 séquence avec un nombre maximum de caractères par ligne
5 séquence avec un nombre maximum de caractères par ligne
6 séquence avec un nombre max
```

La première ligne débute par le caractère > et contient une description de la séquence. Les lignes suivantes contiennent la séquence à proprement dite, mais avec un nombre maximum fixe de caractères par ligne. Ce nombre maximum est habituellement fixé à 60, 70 ou 80 caractères. Une séquence de plusieurs centaines de bases ou de résidus est donc répartie sur plusieurs lignes.

Un fichier est dit *multifasta* lorsqu'il contient plusieurs séquences au format Fasta, les unes à la suite des autres.

Les fichiers contenant une ou plusieurs séquences au format Fasta portent généralement l'extension *.fasta* mais on trouve également *.seq*, *.fas*, *.fna* ou *.faa*.

22.1.1 Exemples

La séquence protéique au format Fasta de la sous-unité β de l'hémoglobine humaine¹, extraite de la base de données UniProt², est :

```
1 >sp|P68871|HBB_HUMAN Hemoglobin subunit beta OS=Homo sapiens OX=9606 GN=HBB PE=1 SV=2
2 MVHLTPEEKSAVTALWGKVVNDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPK
3 VKAHGKKVLGAFSDGLAHLAHLNKGTFATLSELHCDKLHVDPENFRLLGNVLCVLAHFFG
4 KEFTPPVQAAYQKVVAGVANALAHKYH
```

1. <https://www.uniprot.org/uniprot/P68871>

2. <https://www.uniprot.org/>

La première ligne contient la description de la séquence (*Hemoglobin subunit beta*), le type de base de données (ici *sp* qui signifie Swiss-Prot), son identifiant (*P68871*) et son nom (*HBB_HUMAN*) dans cette base de données, ainsi que d'autres informations (*S=Homo sapiens OX=9606 GN=HBB PE=1 SV=2*).

Les lignes suivantes contiennent la séquence sur des lignes ne dépassant pas, ici, 60 caractères. La séquence de la sous-unité β de l'hémoglobine humaine est composée de 147 acides aminés, soit 2 lignes de 40 caractères et une troisième de 27 caractères.

Voici maintenant la séquence nucléique (ARN), au format Fasta, de l'insuline humaine³, extraite de la base de données GenBank⁴ :

```

1 >BT006808.1 Homo sapiens insulin mRNA, complete cds
2 ATGGCCCTGTGGATGCGCCTCCTGCCCTGCTGGCGCTGCTGGCCCTCTGGGGACCTGACCCAGCCGCAG
3 CCTTTGTGAACCAACACCTGTGCGGCTCACACCTGGTGAAGCTCTCTACCTAGTGTGCGGGGAACGAGG
4 CTTCTTACACACCCAAGACCCGCCGGAGGCAGAGGACCTGCAGGTGGGGCAGGTGGAGCTGGGCGGG
5 GGCCCTGGTGCAGGCAGCCTGCAGCCCTTGGCCCTGGAGGGGTCCCTGCAGAAGCGTGGCATTGTGGAAC
6 AATGCTGTACCAGCATCTGCTCCCTTACCAGCTGGAGAACTACTGCAACTAG

```

On retrouve sur la première ligne la description de la séquence (*Homo sapiens insulin mRNA*), ainsi que son identifiant (*BT006808.1*) dans la base de données GenBank.

Les lignes suivantes contiennent les 333 bases de la séquence, réparties sur 5 lignes de 70 caractères maximum.

Pour terminer, voici trois séquences protéiques, au format Fasta, qui correspondent à l'insuline chez humaine (*Homo sapiens*), féline (*Felis catus*) et bovine (*Bos taurus*) :

```

1 >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
2 MALWMRLPLLALLALWGPDPAAAFVNQHLGSHLVEALYLVCGERGFFYTPKTRREAED
3 LQVQVELGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
4 >sp|P06306|INS_FELCA Insulin OS=Felis catus OX=9685 GN=INS PE=1 SV=2
5 MAPWTRLLPLLALLSLWIPAPTRAFVNQHLGSHLVEALYLVCGERGFFYTPKARREAED
6 LQKDAELGEAPGAGGLQPSALEAPLQKRGIVEQCASVCSLYQLEHYCN
7 >sp|P01317|INS_BOVIN Insulin OS=Bos taurus OX=9913 GN=INS PE=1 SV=2
8 MALWTRLRPLLALLALWPPPPARAFVNQHLGSHLVEALYLVCGERGFFYTPKARREVEG
9 PQVGALELAGGPGAGGLEGPPQKRGIVEQCASVCSLYQLENYCN

```

Ces séquences proviennent de la base de données UniProt et sont téléchargeables en suivant ce lien⁵.

Chaque séquence est délimitée par la ligne d'en-tête qui débute par >.

22.1.2 Manipulation avec Python

À partir de l'exemple précédent des 3 séquences d'insuline, voici un exemple de code qui lit un fichier fasta avec Python :

```

1 prot_dict = {}
2 with open("insulin.fasta", "r") as fasta_file:
3     prot_id = ""
4     for line in fasta_file:
5         if line.startswith(">"):

```

3. <https://www.ncbi.nlm.nih.gov/nuccore/BT006808.1?report=fasta>

4. <https://www.ncbi.nlm.nih.gov/nuccore/AY899304.1?report=genbank>

5. <https://www.uniprot.org/uniprot/?sort=score&desc=&compress=no&query=id:P01308%20OR%20id:P01317%20OR%20id:P06306&format=fasta>


```

6     prot_id = line[1:].split()[0]
7     prot_dict[prot_id] = ""
8     else:
9         prot_dict[prot_id] += line.strip()
10    for id in prot_dict:
11        print(id)
12        print(prot_dict[id][:30])

```

Pour chaque séquence lue dans le fichier fasta, on affiche son identifiant et son nom puis les 30 premiers résidus de sa séquence :

```

1 sp|P06306|INS_FELCA
2 MAPWTRLLPLLALLSLWIPAPTRAFVNQHL
3 sp|P01317|INS_BOVIN
4 MALWTRLRPLLALLALWPPPPARAFVNQHL
5 sp|P01308|INS_HUMAN
6 MALWMRLLPLLALLALWGPDPAAAFVNQHL

```

Notez que les protéines sont stockées dans un dictionnaire (`prot_dict`) où les clefs sont les identifiants et les valeurs les séquences.

On peut faire la même chose avec le module BioPython :

```

1 from Bio import SeqIO
2 with open("insulin.fasta", "r") as fasta_file:
3     for record in SeqIO.parse(fasta_file, "fasta"):
4         print(record.id)
5         print(str(record.seq)[:30])

```

Cela produit le même résultat. L'écriture avec BioPython est plus compacte car on utilise ici la fonction `SeqIO.parse()` qui s'occupe de lire le fichier fasta.

22.2 GenBank

GenBank est une banque de séquences nucléiques. Le format de fichier associé contient l'information nécessaire pour décrire un gène ou une portion d'un génome. Les fichiers GenBank portent le plus souvent l'extension `.gbk`.

Le format GenBank est décrit de manière très complète sur le site du NCBI⁶. En voici néanmoins les principaux éléments avec l'exemple du gène qui code pour la trypsine⁷ chez l'Homme.

Remarque : [...] désigne une coupure des auteurs dans le fichier.

22.2.1 L'en-tête

```

1 LOCUS      HUMTRPSGNA                800 bp   mRNA   linear   PRI 14-JAN-1995
2 DEFINITION Human pancreatic trypsin 1 (TRY1) mRNA, complete cds.
3 ACCESSION M22612
4 VERSION   M22612.1
5 KEYWORDS  trypsinogen.

```

6. <https://www.ncbi.nlm.nih.gov/Sitemap/samplerecord.html>

7. <https://www.ncbi.nlm.nih.gov/nuccore/M22612.1>

```

6 SOURCE      Homo sapiens (human)
7   ORGANISM  Homo sapiens
8             Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
9             Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;
10            Catarrhini; Hominidae; Homo.
11 [...]

```

- ligne 1 (LOCUS) : le nom du locus (*HUMTRPSGNA*), la taille du gène (800 paires de base), le type de molécule (ARN messenger)
- ligne 3 (ACCESSION) : l'identifiant de la séquence (*M22612*)
- ligne 4 (VERSION) : la version de la séquence (*M22612.1*). Le nombre qui est séparé de l'identifiant de la séquence par un point est incrémenté pour chaque nouvelle version de la fiche GenBank. Ici *.1* indique que nous en sommes à la première version.
- ligne 6 (SOURCE) : la provenance de la séquence. Souvent l'organisme d'origine.
- ligne 7 (ORGANISME) : le nom scientifique de l'organisme, suivi de sa taxonomie (lignes 8 à 10).

22.2.2 Les features

```

1 [...]
2 FEATURES          Location/Qualifiers
3   source           1..800
4                   /organism="Homo sapiens"
5                   /mol_type="mRNA"
6                   /db_xref="taxon:9606"
7                   /map="7q32-qter"
8                   /tissue_type="pancreas"
9   gene             1..800
10                  /gene="TRY1"
11  CDS               7..750
12                  /gene="TRY1"
13                  /codon_start=1
14                  /product="trypsinogen"
15                  /protein_id="AAA61231.1"
16                  /db_xref="GDB:G00-119-620"
17                  /translation="MNPLLILTFVAAALAAPFDDDDKIVGGYNCEENSVPYQVSLNSG
18                  YHFCGGSLINEQWVVSAGHCYKSRIQVRLGEHNIEVLEGNEQFINAAKIIRHPQYDRK
19                  TLNNDIMLIKLSRAVINARVSTISLPTAPPATGTKCLISGWGNTASSGADYPDELQC
20                  LDAPVLSQAKCEASYPGKITSNMFCVGFLEGGKDSCQGDGGPVCNGQLQGVVSWG
21                  GCAQKNKPGVYTKVYNYVKWIKNTIAANS"
22  sig_peptide       7..51
23                  /gene="TRY1"
24                  /note="G00-119-620"
25 [...]

```

- ligne 9 (gene 1..800) : la délimitation du gène. Ici de la base 1 à la base 800. Par ailleurs, la notation $\langle x..y \rangle$ indique que la séquence est partielle sur l'extrémité 5'. Réciproquement, $\langle x..y \rangle$ indique que la séquence est partielle sur l'extrémité 3'. Enfin, pour les séquences d'ADN, la notation `complement(x..y)` indique que le gène se trouve de la base x à la base y , mais sur le brin complémentaire.
- ligne 10 (`/gene="TRY1"`) : le nom du gène
- ligne 11 (CDS 7..750) : la délimitation de la séquence codante
- ligne 14 (`/product="trypsinogen"`) : le nom de la protéine produite

- ligne 17 à 20 (`/translation="MNPLLLIL..."`) : la séquence protéique issue de la traduction de la séquence codante
- ligne 22 (`sig_peptide 7..51`) : la délimitation du peptide signal

22.2.3 La séquence

```

1 [...]
2 ORIGIN
3     1 accaccatga atccactcct gatccttacc tttgtggcag ctgctcttgc tgcccccttt
4     61 gatgatgatg acaagatcgt tgggggctac aactgtgagg agaattctgt cccctaccag
5    121 gtgtccctga attctggcta ccacttctgt ggtggctccc tcatcaacga acagtgggtg
6    181 gtatcagcag gccactgcta caagtccgc atccaggtga gactgggaga gcacaacatc
7    241 gaagtcctgg aggggaatga gcagttcatc aatgcagcca agatcatccg ccaccccaaa
8    301 tacgacagga agactctgaa caatgacatc atgttaatca agctctcctc acgtgcagta
9    361 atcaacgccg gcgtgtccac catctctctg cccaccgccg ctccagccac tggcacgaag
10   421 tgcctcatct ctggctgggg caaactgctg agctctggcg cggactacc agacgagctg
11   481 cagtgcttgg atgctcctgt gctgagccag gctaagtgtg aagcctccta ccctggaaag
12   541 attaccagca acatgttctg tgtgggcttc cttgaggag gcaaggattc atgtcagggt
13   601 gattctggtg gccctgtggt ctgcaatgga cagctccaag gagttgtctc ctggggtgat
14   661 ggctgtgccc agaagaacaa gcctggagtc tacaccaagg tctacaacta cgtgaaatgg
15   721 attaagaaca ccatagctgc caatagctaa agccccagc atctcttcag tctctatacc
16   781 aataaagtga ccctgttctc
17 //

```

La séquence est contenue entre les balises `ORIGIN` (ligne 2) et `//` (ligne 17).

Chaque ligne est composée d'une série d'espaces, puis du numéro du premier nucléotide de la ligne, puis d'au plus 6 blocs de 10 nucléotides. Chaque bloc est précédé d'un espace.

Par exemple, ligne 10, le premier nucléotide de la ligne (`t`) est le numéro 421 dans la séquence.

22.2.4 Manipulation avec Python

À partir de l'exemple précédent, voici comment lire un fichier GenBank avec Python et le module `BioPython` :

```

1 from Bio import SeqIO
2 with open("M22612.gbk", "r") as gbk_file:
3     record = SeqIO.read(gbk_file, "genbank")
4     print(record.id)
5     print(record.description)
6     print(record.seq[:60])

```

Pour séquence lue dans le fichier GenBank, on affiche son identifiant, sa description et les 60 premiers résidus :

```

1 M22612.1
2 Human pancreatic trypsin 1 (TRY1) mRNA, complete cds.
3 ACCACCATGAATCCACTCCTGATCCTTACCTTTGTGGCAGCTGCTCTTGCTGCCCTTT

```

Il est également possible de lire un fichier GenBank sans le module `BioPython`. Une activité dédiée est proposée dans le chapitre *Mini-projets*.

22.3 PDB

La Protein Data Bank⁸ (PDB) est une banque de données qui contient les structures de biomacromolécules (protéines, ADN, ARN, virus...). Historiquement, le format de fichier qui y est associé est le PDB, dont une documentation détaillée est disponible sur le site éponyme⁹. Les principales extensions de ce format de fichier sont `.ent` et surtout `.pdb`.

Un fichier PDB est constitué de deux parties principales : l'en-tête et les coordonnées. L'en-tête est lisible et utilisable par un être humain (et parfois par une machine). À l'inverse les coordonnées sont surtout utilisables par une programme pour calculer certaines propriétés de la structure ou simplement la représenter pour permettre sa visualisation, par exemple sur l'écran d'un ordinateur. Bien sur, un utilisateur expérimenté peut parfaitement jeter un oeil à cette seconde partie.

Examinons ces deux parties avec la trypsine bovine¹⁰.

22.3.1 En-tête

Pour la trypsine bovine, l'en-tête compte 510 lignes. En voici quelques unes :

```

1  HEADER   HYDROLASE (SERINE PROTEINASE)      26-OCT-81  2PTN
2  TITLE    ON THE DISORDERED ACTIVATION DOMAIN IN TRYPSINOGEN.
3  TITLE    2 CHEMICAL LABELLING AND LOW-TEMPERATURE CRYSTALLOGRAPHY
4  COMPND   MOL_ID: 1;
5  COMPND   2 MOLECULE: TRYPSIN;
6  COMPND   3 CHAIN: A;
7  [...]
8  SOURCE   2 ORGANISM_SCIENTIFIC: BOS TAURUS;
9  [...]
10 EXPDTA   X-RAY DIFFRACTION
11 [...]
12 REMARK   2 RESOLUTION.   1.55 ANGSTROMS.
13 [...]
14 DBREF    2PTN A  16  245 UNP   P00760  TRY1_BOVIN   21   243
15 SEQRES   1 A  223 ILE VAL GLY GLY TYR THR CYS GLY ALA ASN THR VAL PRO
16 SEQRES   2 A  223 TYR GLN VAL SER LEU ASN SER GLY TYR HIS PHE CYS GLY
17 SEQRES   3 A  223 GLY SER LEU ILE ASN SER GLN TRP VAL VAL SER ALA ALA
18 SEQRES   4 A  223 HIS CYS TYR LYS SER GLY ILE GLN VAL ARG LEU GLY GLU
19 [...]
20 HELIX    1  H1 SER A 164 ILE A 176 1SNGL ALPHA TURN,REST IRREG.    13
21 HELIX    2  H2 LYS A 230 VAL A 235 5CONTIGUOUS WITH H3             6
22 HELIX    3  H3 SER A 236 ASN A 245 1CONTIGUOUS WITH H2            10
23 SHEET    1   A 7 TYR A 20  THR A 21  0
24 SHEET    2   A 7 LYS A 156 PRO A 161 -1 N CYS A 157  0  TYR A 20
25 [...]
26 SSBOND   1 CYS A  22   CYS A 157                1555 1555 2.04
27 SSBOND   2 CYS A  42   CYS A  58                1555 1555 2.02
28 [...]

```

— La ligne `HEADER` contient le nom de la protéine (*HYDROLASE (SERINE PROTEINASE)*), la date de dépôt de cette structure dans la banque de données (26 octobre 1981) et l'identifiant de

8. <https://www.rcsb.org/>

9. <http://www.wwpdb.org/documentation/file-format-content/format33/v3.3.html>

10. <https://www.rcsb.org/structure/2PTN>

la structure dans la PDB, on parle parfois de *code PDB* (2PTN).

- Les lignes `TITLE` correspondent au titre de l'article scientifique dans lequel a été publié cette structure.
- Les lignes `COMPND` indiquent que la trypsine est composé d'une seule chaîne peptidique, appelée ici A.
- La ligne `SOURCE` indique le nom scientifique de l'organisme dont provient cette protéine (ici, la boeuf).
- La ligne `EXPDTA` précise la technique expérimentale employée pour déterminer cette structure. Ici, la cristallographie aux rayons X. Mais on peut également trouver *SOLUTION NMR* pour de la résonance magnétique nucléaire en solution, *ELECTRON MICROSCOPY* pour de la microscopie électronique...
- La ligne `REMARK 2` précise, dans le cas d'une détermination par cristallographie aux rayons X, la résolution obtenue, ici 1,55 Angstroms.
- La ligne `DBREF` indique les liens éventuels vers d'autres banques de données. Ici, l'identifiant correspondant à cette protéine dans UniProt (*UNP*) est P00760¹¹.
- Les lignes `SEQRES` donnent à la séquence de la protéine. Les résidus sont représentés par leur code à 3 lettres.
- Les lignes `HELIX` et `SHEET` correspondent aux structures secondaires hélices α et brin β de cette protéine. Ici, *H1 SER A 164 ILE A 176* indique qu'il y a une première hélice α (*H1*) comprise entre les résidus Ser164 et Ile176 de la chaîne A.
- Enfin, les lignes `SSBOND` définissent les bonds disulfures. Par exemple ici entre les résidus Cys22 et Cys157 et entre les résidus Cys42 et Cys58.

22.3.2 Coordonnées

Avec la même protéine, la partie coordonnées représente plus de 1700 lignes. En voici quelques unes correspondantes au résidu leucine 99 :

```

1 [...]
2 ATOM 601 N LEU A 99 10.007 19.687 17.536 1.00 12.25 N
3 ATOM 602 CA LEU A 99 9.599 18.429 18.188 1.00 12.25 C
4 ATOM 603 C LEU A 99 10.565 17.281 17.914 1.00 12.25 C
5 ATOM 604 O LEU A 99 10.256 16.101 18.215 1.00 12.25 O
6 ATOM 605 CB LEU A 99 8.149 18.040 17.853 1.00 12.25 C
7 ATOM 606 CG LEU A 99 7.125 19.029 18.438 1.00 18.18 C
8 ATOM 607 CD1 LEU A 99 5.695 18.554 18.168 1.00 18.18 C
9 ATOM 608 CD2 LEU A 99 7.323 19.236 19.952 1.00 18.18 C
10 [...]
```

Chaque ligne correspond à un atome et débute par `ATOM` ou `HETATM`. `ATOM` désigne un atome de la structure de la biomolécules. `HETATM` est utilisé pour les atomes qui ne sont pas dans une biomolécules, comme les ions ou les molécules d'eau.

Toutes les lignes de coordonnées ont sensiblement le même format. Par exemple pour la première ligne :

- `ATOM` (ou `HETATM`)
- `601` : le numéro de l'atome.
- `N` : le nom de l'atome. Ici, un atome d'azote du squelette peptidique. La structure complète du résidu leucine est représentée figure 22.1
- `Leu` : le résidu dont fait partie l'atome. Ici une leucine.
- `A` : le nom de la chaîne peptidique.
- `99` : le numéro du résidu dans la protéine.

11. <https://www.uniprot.org/uniprot/P00760>

- 10.007 : la coordonnées x de l'atome.
- 19.687 : la coordonnées y de l'atome.
- 17.536 : la coordonnées z de l'atome.
- 1.00 : la probabilité de trouver l'atome à cette position dans l'espace. Cette probabilité est différente de 1 lorsque la technique expérimentale n'a pas déterminé avec une totale certitude la position de l'atome. Par exemple dans le cas d'un atome très mobile dans une structure, qui est déterminé comme étant à deux positions possibles, chaque position aura alors la probabilité 0,50.
- 12.25 : le facteur de température qui est proportionnel à la mobilité de l'atome dans l'espace. Les atomes situés en périphérie d'une structure sont souvent plus mobiles que ceux situés au coeur de la structure.
- N : l'élément chimique de l'atome. Ici, l'azote.

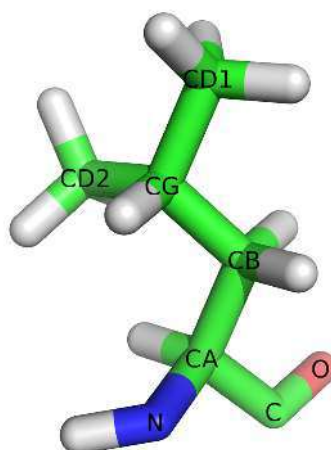


FIGURE 22.1 – Structure tridimensionnelle d'un résidu leucine. Les noms des atomes sont indiqués en noir.

Une documentation plus complète des différents champs qui constituent une ligne de coordonnées atomiques se trouve sur le site de la PDB ¹².

Les résidus sont ensuite décrits les uns après les autres, atome par atome. Voici par exemple les premiers résidus de la trypsine bovine :

```

1 [...]
2 ATOM      1  N   ILE A 16   -8.155  9.648 20.365 1.00 10.68      N
3 ATOM      2  CA  ILE A 16   -8.150  8.766 19.179 1.00 10.68      C
4 ATOM      3  C   ILE A 16   -9.405  9.018 18.348 1.00 10.68      C
5 ATOM      4  O   ILE A 16  -10.533  8.888 18.870 1.00 10.68      O
6 ATOM      5  CB  ILE A 16   -8.091  7.261 19.602 1.00 10.68      C
7 ATOM      6  CG1 ILE A 16   -6.898  6.882 20.508 1.00  7.42      C
8 ATOM      7  CG2 ILE A 16   -8.178  6.281 18.408 1.00  7.42      C
9 ATOM      8  CD1 ILE A 16   -5.555  6.893 19.773 1.00  7.42      C
10 ATOM     9  N   VAL A 17   -9.224  9.305 17.090 1.00  9.63      N
11 ATOM     10  CA  VAL A 17  -10.351  9.448 16.157 1.00  9.63      C
12 ATOM     11  C   VAL A 17  -10.500  8.184 15.315 1.00  9.63      C
13 ATOM     12  O   VAL A 17   -9.496  7.688 14.748 1.00  9.63      O
14 ATOM     13  CB  VAL A 17  -10.123 10.665 15.222 1.00  9.63      C
15 ATOM     14  CG1 VAL A 17  -11.319 10.915 14.278 1.00 11.95      C
16 ATOM     15  CG2 VAL A 17   -9.737 11.970 15.970 1.00 11.95      C
17 [...]

```

12. <http://www ww p d b . o r g / d o c u m e n t a t i o n / f i l e - f o r m a t - c o n t e n t / f o r m a t 3 3 / s e c t 9 . h t m l>

Vous remarquerez que le numéro du premier résidu est 16 et non pas. Cela s'explique par la technique expérimentale utilisée qui n'a pas permis de déterminer la structure des 15 premiers résidus.

La structure de la trypsine bovine n'est constituée que d'une seule chaîne (notée A). Lorsqu'une structure est composée de plusieurs chaînes, comme dans le cas de la structure du récepteur GABAB 1 et 2 chez la drosophile (code PDB 5X9X¹³) :

```

1 [...]
2 ATOM 762 HB1 ALA A 44 37.162 -2.955 2.220 1.00 0.00 H
3 ATOM 763 HB2 ALA A 44 38.306 -2.353 3.417 1.00 0.00 H
4 ATOM 764 HB3 ALA A 44 38.243 -1.621 1.814 1.00 0.00 H
5 TER 765 ALA A 44
6 ATOM 766 N GLY B 95 -18.564 3.009 13.772 1.00 0.00 N
7 ATOM 767 CA GLY B 95 -19.166 3.646 12.621 1.00 0.00 C
8 ATOM 768 C GLY B 95 -20.207 2.755 11.976 1.00 0.00 C
9 [...]
```

La première chaîne est notée A et la seconde B. La séparation entre les deux est marquée par la ligne
 TER 765 ALA A 44.

Enfin, lorsque la structure est déterminée par RMN, il est possible que plusieurs structures soient présentes dans le même fichier PDB. Toutes ces structures, ou modèles, sont des solutions possibles du jeu de contraintes mesurées expérimentalement en RMN. Voici un exemple, toujours pour structure du récepteur GABAB 1 et 2 chez la drosophile :

```

1 [...]
2 MODEL 1
3 ATOM 1 N MET A 1 -27.283 -9.772 5.388 1.00 0.00 N
4 ATOM 2 CA MET A 1 -28.233 -8.680 5.682 1.00 0.00 C
5 [...]
6 ATOM 1499 HG2 GLU B 139 36.113 -5.242 2.536 1.00 0.00 H
7 ATOM 1500 HG3 GLU B 139 37.475 -4.132 2.428 1.00 0.00 H
8 TER 1501 GLU B 139
9 ENDMDL
10 MODEL 2
11 ATOM 1 N MET A 1 -29.736 -10.759 4.394 1.00 0.00 N
12 ATOM 2 CA MET A 1 -28.372 -10.225 4.603 1.00 0.00 C
13 [...]
14 ATOM 1499 HG2 GLU B 139 36.113 -5.242 2.536 1.00 0.00 H
15 ATOM 1500 HG3 GLU B 139 37.475 -4.132 2.428 1.00 0.00 H
16 TER 1501 GLU B 139
17 ENDMDL
18 MODEL 2
19 ATOM 1 N MET A 1 -29.736 -10.759 4.394 1.00 0.00 N
20 ATOM 2 CA MET A 1 -28.372 -10.225 4.603 1.00 0.00 C
21 [...]
```

Chaque structure est encadrée par les lignes

```
1 MODEL n
```

et

```
1 ENDMDL
```

13. <http://www.rcsb.org/structure/5X9X>

où n est le numéro du modèle. Pour la structure du récepteur GABAB 1 et 2, il y a 20 modèles de décrits dans le fichier PDB.

22.3.3 Manipulation avec Python

Le module BioPython peut également lire un fichier PDB.

Chargement de la structure de la trypsine bovine :

```
1 from Bio.PDB import PDBParser
2 parser = PDBParser()
3 prot_id = "2PTN"
4 prot_file = "2PTN.pdb"
5 structure = parser.get_structure(prot_id, prot_file)
```

Remarque : les fichiers PDB sont parfois (très) mal formatés. Si BioPython ne parvient pas à lire un tel fichier, remplacez alors la 2e ligne par `parser = PDBParser(PERMISSIVE=1)`. Soyez néanmoins très prudent quant aux résultats obtenus.

Affichage du nom de la structure et de la technique expérimentale utilisée pour déterminer la structure :

```
1 print(structure.header["head"])
2 print(structure.header["structure_method"])
```

ce qui produit :

```
1 hydrolase (serine proteinase)
2 x-ray diffraction
```

Extraction des coordonnées de l'atome N du résidu Ile16 et de l'atome CA du résidu Val17 :

```
1 model = structure[0]
2 chain = model["A"]
3 res1 = chain[16]
4 res2 = chain[17]
5 print(res1.resname, res1["N"].coord)
6 print(res2.resname, res2["CA"].coord)
```

ce qui produit :

```
1 ILE [ -8.15499973 9.64799976 20.36499977]
2 VAL [-10.35099983 9.44799995 16.15699959]
```

Remarque : BioPython utilise la hiérarchie suivante : `structure > model > chain > residue > atom`, même lorsque la structure ne contient qu'un seul modèle, comme c'est le cas ici, puisque la structure a été obtenue par cristallographie aux rayons X.

Enfin, pour afficher les coordonnées des carbones α (notés CA) des 10 premiers résidus (à partir du résidu 16 car c'est le premier résidu dont on connaît la structure) :

```
1 res_start = 16
2 model = structure[0]
3 chain = model["A"]
4 for i in range(10):
5     idx = res_start + i
6     print(chain[idx].resname, idx, chain[idx]["CA"].coord)
```


avec pour résultat :

```

1 ILE 16 [ -8.14999962 8.76599979 19.17900085]
2 VAL 17 [-10.35099983 9.44799995 16.15699959]
3 GLY 18 [-12.02099991 6.63000011 14.25899982]
4 GLY 19 [-10.90200043 3.89899993 16.68400002]
5 TYR 20 [-12.65100002 1.44200003 19.01600075]
6 THR 21 [-13.01799965 0.93800002 22.76000023]
7 CYS 22 [-10.02000046 -1.16299999 23.76000023]
8 GLY 23 [-11.68299961 -2.86500001 26.7140007 ]
9 ALA 24 [-10.64799976 -2.62700009 30.36100006]
10 ASN 25 [ -6.96999979 -3.43700004 31.02000046]

```

Il est aussi très intéressant (et formateur) d'écrire son propre *parser* de fichier PDB, c'est-à-dire un programme qui lit un fichier PDB (sans le module BioPython). Dans ce cas, la figure 22.2 vous aidera à déterminer comment extraire les différentes informations d'une ligne de coordonnées ATOM ou HETATM.

Exemple : pour extraire le nom du résidu, il faut isoler le contenu des colonnes 18 à 20 du fichier PDB, ce qui correspond aux index de 17 à 19 pour une chaîne de caractère en Python, soit la tranche de chaîne de caractères [17:20] car la première borne est incluse et la seconde exclus.

Pour lire le fichier PDB de trypsine bovine (2PTN.pdb) et extraire (encore) les coordonnées des carbones α des 10 premiers résidus, nous pouvons utiliser le code suivant :

```

1 with open("2PTN.pdb", "r") as pdb_file:
2     res_count = 0
3     for line in pdb_file:
4         if line.startswith("ATOM"):
5             atom_name = line[12:16].strip()
6             res_name = line[17:20].strip()
7             res_num = int(line[22:26])
8             if atom_name == "CA":
9                 res_count += 1
10                x = float(line[30:38])
11                y = float(line[38:46])
12                z = float(line[46:54])
13                print(res_name, res_num, x, y, z)
14            if res_count >= 10:
15                break

```

ce qui donne :

```

1 ILE 16 -8.15 8.766 19.179
2 VAL 17 -10.351 9.448 16.157
3 GLY 18 -12.021 6.63 14.259
4 GLY 19 -10.902 3.899 16.684
5 TYR 20 -12.651 1.442 19.016
6 THR 21 -13.018 0.938 22.76
7 CYS 22 -10.02 -1.163 23.76
8 GLY 23 -11.683 -2.865 26.714
9 ALA 24 -10.648 -2.627 30.361
10 ASN 25 -6.97 -3.437 31.02

```

Remarque : pour extraire des valeurs numériques, comme des numéros de résidus ou des coordonnées atomiques, il ne faudra pas oublier de les convertir en entiers ou en réels.

