

# **Bioinformatique avec Python et Biopython**



# Bioinformatique avec Python et Biopython

**Emmanuel Jaspard**  
Université d'Angers

**Gilles Hunault**  
Université d'Angers

**DUNOD**

#### **NOUS NOUS ENGAGEONS EN FAVEUR DE L'ENVIRONNEMENT :**



Nos livres sont imprimés sur des papiers certifiés pour réduire notre impact sur l'environnement.



Le format de nos ouvrages est pensé afin d'optimiser l'utilisation du papier.



Depuis plus de 30 ans, nous imprimons 70 % de nos livres en France et 25 % en Europe et nous mettons tout en œuvre pour augmenter cet engagement auprès des imprimeurs français.



Nous limitons l'utilisation du plastique sur nos ouvrages (film sur les couvertures et les livres).

© Dunod, Paris, 2024  
11, rue Paul Bert, 92240 Malakoff  
[www.dunod.com](http://www.dunod.com)  
ISBN 978-2-10-085887-3

À Gwilherm Jaspard.

Pour Raphaël et Timothée.





# Table des matières

Introduction	1
<b>Chapitre 1</b> Dénombrement des nucléotides et calculs associés (%GC, $T_m$ )	7
1. Rappel sur la structure des nucléotides	7
2. La composition en nucléotides de l'ADN et des ARN	8
3. Illustration des îlots CpG	9
4. Calculs du pourcentage de GC et de la température de fusion $T_m$	10
5. Buts du script <code>gctm.py</code>	11
6. Code du script <code>gctm.py</code>	12
7. Explications du code <code>gctm.py</code>	12
8. Code de l'auto-test pour le script <code>gctm.py</code>	14
9. Calculs de %GC et de $T_m$ avec Biopython	15
L'essentiel	16
Bibliographie	16
Ressources Web	16
Exercices	18
<b>Chapitre 2</b> Écriture de séquences d'acides nucléiques	19
1. Les séquences d'acides nucléiques	19
2. Les formats des fichiers bioinformatiques – le format <b>FASTA</b>	20
3. Buts du script <code>brinscomp.py</code>	22
4. Code du script <code>brinscomp.py</code>	23
5. Explications du code <code>brinscomp.py</code>	23
6. Code de l'auto-test pour le script <code>brinscomp.py</code>	24
L'essentiel	25
Bibliographie	25

Ressources Web	25
Exercices	26
<b>Chapitre 3 Séquences et expressions régulières</b>	<b>27</b>
1. Rappels sur les caractéristiques des expressions régulières	27
2. Utilisation des expressions régulières dans l'analyse des données bioinformatiques	28
3. Buts du script <code>rchmotif.py</code>	29
4. Code du script <code>rchmotif.py</code>	30
5. Explications du code <code>rchmotif.py</code>	31
6. Code de l'auto-test pour le script <code>rchmotif.py</code>	31
L'essentiel	32
Bibliographie	32
Ressources Web	32
Exercices	33
<b>Chapitre 4 Recherches bibliographiques dans PubMed</b>	<b>35</b>
1. Rappel sur les données bibliographiques et le moteur de recherche <code>Entrez</code>	35
2. Buts du script <code>rchbiblio.py</code>	37
3. Code du script <code>rchbiblio.py</code>	38
4. Explications du code <code>rchbiblio.py</code>	39
5. Code de l'auto-test pour le script <code>rchbiblio.py</code>	40
L'essentiel	41
Bibliographie	41
Ressources Web	41
Exercices	42

<b>Chapitre 5</b>	<b>Profil d'une famille de protéines</b>	43
	1. Rappel sur les domaines structuraux des protéines	43
	2. Profils, domaines et bases de données de familles de protéines	44
	3. Méthode de construction des alignements <b>PFAM</b>	45
	4. Buts du script <code>familprot.py</code>	45
	5. Code du script <code>familprot.py</code>	47
	6. Explications du code <code>familprot.py</code>	48
	7. Code de l'auto-test pour le script <code>familprot.py</code>	51
	L'essentiel	52
	Bibliographie	52
	Ressources Web	52
	Exercices	53
<b>Chapitre 6</b>	<b>Profil d'hydrophobicité d'une protéine – fenêtre de taille variable</b>	55
	1. Propriétés physico-chimiques des acides aminés	55
	2. Calcul de l'hydrophobicité d'une séquence d'acides aminés avec une fenêtre glissante	56
	3. Exemple d'un récepteur couplé à une protéine G ( <b>RCPG</b> )	57
	4. Buts du script <code>profilhyd.py</code>	57
	5. Code du script <code>profilhyd.py</code>	58
	6. Explications du code <code>profilhyd.py</code>	59
	7. Code de l'auto-test pour le script <code>profilhyd.py</code>	60
	8. Code du script <code>profilhyd_utils.py</code>	61
	9. Explications du code <code>profilhyd_utils.py</code>	62
	10. Calculs avec les fonctions de <b>Biopython</b>	62
	L'essentiel	63
	Bibliographie	63

Ressources Web	63
Exercices	64
<b>Chapitre 7</b> <b>Position de cystéines impliquées dans un pont disulfure</b>	65
1. Propriétés physico-chimiques et particularité de la cystéine	65
2. Pont disulfure et structure 3D des protéines	66
3. La base de données « <i>Protein DataBank</i> » ou PDB	67
4. Buts du script <code>disulfure.py</code>	69
5. Code du script <code>disulfure.py</code>	70
6. Explications du code <code>disulfure.py</code>	71
7. Code de l'auto-test pour le script <code>disulfure.py</code>	73
L'essentiel	74
Bibliographie	74
Ressources Web	74
Exercices	75
<b>Chapitre 8</b> <b>Calcul du barycentre d'une chaîne polypeptidique repliée</b>	77
1. Structure tridimensionnelle des protéines	77
2. Calcul du barycentre des carbones $\alpha$ des acides aminés d'une chaîne polypeptidique repliée	78
3. Buts du script <code>barycentre.py</code>	78
4. Code du script <code>barycentre.py</code>	79
5. Explications du code <code>barycentre.py</code>	80
6. Code de l'auto-test pour le script <code>barycentre.py</code>	82
L'essentiel	83
Bibliographie	83
Ressources Web	83
Exercices	84

<b>Chapitre 9</b>	<b>Diagrammes de Ramachandran</b>	85
	1. La liaison peptidique	85
	2. Les angles de torsion	86
	3. Les diagrammes de Ramachandran	87
	4. Buts du script <code>ramachandran.py</code>	87
	5. Code du script <code>ramachandran.py</code>	89
	6. Explications du code <code>ramachandran.py</code>	90
	7. Code de l'auto-test pour le script <code>ramachandran.py</code>	91
	8. Code du script <code>ramachandran_utils.py</code>	92
	9. Explications du code <code>ramachandran_utils.py</code>	94
	L'essentiel	96
	Bibliographie	96
	Ressources Web	96
	Exercices	97
<b>Chapitre 10</b>	<b>Contacts entre résidus d'acides aminés dans les structures 3D</b>	99
	1. Interactions entre résidus d'acides aminés au sein des protéines	99
	2. Contacts et distances entre résidus d'acides aminés	100
	3. Buts du script <code>contacts.py</code>	101
	4. Code du script <code>contacts.py</code>	102
	5. Explications du code <code>contacts.py</code>	102
	6. Code de l'auto-test pour le script <code>contacts.py</code>	103
	7. Code du script <code>structure_pdb.py</code>	103
	8. Explications du code <code>structure_pdb.py</code>	104
	9. Code du script <code>contacts_utils.py</code>	105
	10. Explications du code <code>contacts_utils.py</code>	106

L'essentiel	108
Bibliographie	108
Ressources Web	108
Exercices	109
<b>Chapitre 11 Superposition de structures protéiques</b>	111
1. Superposition de structures de protéines et scores de distances interatomiques	111
2. Exemples de scores de mesure de similarité de structures	112
3. L'algorithme <code>AlphaFold</code>	113
4. Buts du script <code>superp.py</code>	114
5. Code du script <code>superp.py</code>	115
6. Explications du code <code>superp.py</code>	116
7. Code de l'auto-test pour le script <code>superp.py</code>	116
8. Code du script <code>superp_pdb.py</code>	117
9. Explications du code <code>superp_pdb.py</code>	118
L'essentiel	120
Bibliographie	120
Ressources Web	120
Exercices	121
<b>Chapitre 12 Les protéines ou régions intrinsèquement désordonnées</b>	123
1. Découverte des protéines ou régions intrinsèquement désordonnées	123
2. Abondance et fonctions biologiques des <code>IDP/IDR</code>	125
3. Prédiction des <code>IDP/IDR</code>	126
4. Buts du script <code>repliement.py</code>	126
5. Code du script <code>repliement.py</code>	127

6. Explications du code <code>repliement.py</code>	128
7. Code de l'auto-test pour le script <code>repliement.py</code>	129
L'essentiel	130
Bibliographie	130
Ressources Web	130
Exercices	131
<b>Chapitre 13</b> <b>Alignements de séquences et protéines homologues</b>	133
1. Notions élémentaires d'alignement de séquences	133
2. Description de <b>BLAST</b> (« <i>Basic Local Alignment Search Tool</i> »)	135
3. Principe de <b>PHI-BLAST</b> (« <i>Pattern Hit Initiated-BLAST</i> »)	136
4. Buts du script <code>homologues.py</code>	138
5. Code du script <code>homologues.py</code>	139
6. Explications du code <code>homologues.py</code>	140
7. Code de l'auto-test pour le script <code>homologues.py</code>	141
8. Code du script <code>homologues_utils.py</code>	142
9. Explications du code <code>homologues_utils.py</code>	144
L'essentiel	146
Bibliographie	146
Ressources Web	146
Exercices	147
<b>Réponses aux questions</b>	149
Chapitre 1	149
Chapitre 2	152
Chapitre 3	155
Chapitre 4	158

## Table des matières

Chapitre 5	160
Chapitre 6	161
Chapitre 7	165
Chapitre 8	166
Chapitre 9	167
Chapitre 10	169
Chapitre 11	171
Chapitre 12	172
Chapitre 13	176
<b>Annexes</b>	181
<b>Annexe 1 : rappels Python</b>	183
<b>Annexe 2 : expressions régulières</b>	191
<b>Annexe 3 : présentation de Biopython</b>	197
<b>Annexe 4 : liste des scripts dans les chapitres</b>	199
<b>Annexe 5 : liste des scripts dans les réponses aux questions</b>	201
<b>Annexe 6 : packages et modules cités</b>	203
<b>Index</b>	205

# Introduction

La bioinformatique est une discipline essentielle de la biologie. Il existe différentes acceptions du terme « bioinformatique » : elle peut être définie comme l'ensemble des théories et des méthodes issues de plusieurs disciplines, notamment l'informatique, dédiées aux traitements des données biologiques.

Les domaines dits en « omique » (dont la génomique, la métagénomique, la transcriptomique, la protéomique, la métabolomique, l'interactomique, la fluxomique...) nécessitent des outils bioinformatiques. Compte tenu de la très grande quantité de données biologiques acquises à une vitesse croissante et en regard de leur complexité, notamment la diversité des types de données, il est impossible d'en effectuer une analyse globale et exhaustive sans programme informatique.

Certains langages de programmation sont donc des outils indispensables pour l'acquisition des données biologiques et leur intégration dans des bases de données, la transformation des types de données, la fouille de données et leur classification, l'analyse (numérique, statistique, textuelle, visuelle...) des données selon des modèles mathématiques ou physiques et bien d'autres types de traitements.

Nous avons choisi `Python` et `Biopython` pour les raisons suivantes :

- Le langage `Python` est un langage orienté objet particulièrement utilisé en bioinformatique : c'est un langage évolué relativement simple à apprendre, qui est adopté par une proportion croissante de chercheurs pour développer des applications en bioinformatique, avec une multiplicité de bibliothèques et de scripts libres de droits...
- `Biopython` est un ensemble d'outils, librement disponibles et issus d'un effort collaboratif, pour développer des bibliothèques et des applications `Python` pour les besoins des travaux en bioinformatique.

La première partie de cet ouvrage a trait à l'analyse de certaines propriétés et caractéristiques des chaînes de nucléotides (**ADN** et **ARN**). La seconde partie traite de l'analyse textuelle (fouille de données de type chaînes de caractères) avec un outil remarquable, les expressions régulières. La troisième partie est dédiée à l'analyse des chaînes polypeptidiques et aux caractéristiques physico-chimiques des chaînes latérales des acides aminés. Elle illustre aussi l'alignement de séquences en utilisant la puissance des expressions régulières dans l'outil **PHI-BLAST**.

Nous explicitons dans cet ouvrage certaines interfaces de programmation `Biopython` (**API**) afin que chacun comprenne comment communiquer avec telle ou telle base de données et avec les interfaces *Web* des grands sites internationaux de bioinformatique.

Bon nombre de scripts équivalents aux scripts de cet ouvrage existent sous de multiples formes et sont accessibles librement dans des dépôts tels que **GitHub**. Même si nos scripts sont complets, fonctionnels et téléchargeables, le but principal de cet ouvrage n'est pas de fournir des scripts « clé en main », mais bien d'en comprendre la logique et

d'établir un lien avec le concept biologique afférent afin que le lecteur puisse développer ses propres scripts ou adapter ceux existants. Pour continuer l'analogie automobile, nous fournissons des scripts détaillés et expliqués pour que les lecteurs et lectrices aillent « voir sous le capot ».

## Contenu type d'un chapitre

---

Tous les chapitres sont structurés de la même façon : on trouve d'abord un rappel sur les notions biologiques liées au sujet du chapitre. Ensuite, les objectifs du script sont présentés, avec un exemple de résultat de son exécution. Le code du script est alors fourni, accompagné d'explications et du code de l'auto-test utilisé pour vérifier les principaux cas d'exécution. Lorsque le script est un peu long, les sous-programmes (fonctions) développés pour le script sont listés et expliqués. Enfin, une synthèse du chapitre précède une courte bibliographie et une liste de sites et pages *Web* susceptibles d'éclairer ou d'approfondir le chapitre. En fin de chapitre, des questions complètent aussi bien les aspects biologiques que le code `Python` utilisé. Les réponses détaillées à ces questions sont disponibles après le chapitre 13.

Aucun des scripts ne fournit d'interface de saisie utilisateur, excepté au chapitre 1. Les scripts implémentent toujours des fonctions paramétrées. Le chapitre 1 décrit comment réaliser une interface minimale avec une lecture au clavier, sachant que ce n'est pas toujours la meilleure solution. Une interface *Web*, par exemple, ou une lecture dans un fichier (local ou sur Internet) sont parfois plus adaptées. C'est pourquoi nous fournissons dans chaque chapitre l'auto-test lié au script afin de savoir comment interfacier (au sens d'utiliser) les fonctions implémentées.

## Progression pédagogique : du plus simple au plus compliqué

---

Dans le chapitre 1, nous nous intéressons à des calculs simples : dénombrement de nucléotides, calcul du pourcentage de GC, de la température de fusion  $T_m$ . Nous effectuons directement ces calculs en `Python` avant d'expliquer comment utiliser les fonctions équivalentes de `Biopython`. Nous montrons aussi comment profiter de l'auto-test *via* le test `if __name__=="__main__"` afin de tester et d'exposer le comportement des fonctions.

Le chapitre 2 traite de la lecture et de l'écriture de séquences d'acides nucléiques (**ADN** ou **ARN**) d'une séquence lue dans un fichier local.

Le chapitre 3 est consacré aux expressions régulières et à la recherche de motifs dans des fichiers au format **FASTA**.

Le chapitre 4 présente une utilisation typique de `Biopython` : on utilise les méthodes fournies par le package **Entrez** afin de trouver des articles dans la base de données bibliographiques du site **PubMed**. La connexion au serveur *Web* du **NCBI** et

la gestion des fichiers **XML** renvoyés y sont « transparentes » pour le développeur parce que le package implémente des structures de données et des méthodes qui les gèrent de façon automatique.

Le chapitre 5 présente les profils de familles de protéines et la recherche de familles **PFAM** sur le site **EBI-Interpro**. Le script commence par rechercher le numéro identifiant d'une famille de protéines à partir d'un mot-clé avant de rapatrier le profil associé.

Le chapitre 6 détaille le profil d'hydrophobicité d'une protéine dans une fenêtre glissante de taille choisie. Ce calcul est effectué une première fois en **Python** « pur », avant d'être réalisé par l'appel d'une fonction du package **Bio.SeqUtils.ProtParam** de **Biopython**. Les fichiers utilisés, au format **FASTA**, sont soit locaux, soit lus sur le site **Uniprot**.

Le chapitre 7 expose la recherche de la position de cystéines impliquées dans un ou plusieurs ponts disulfures pour un fichier structure d'une protéine lu dans la **PDB**.

Le chapitre 8 explicite le calcul du barycentre des carbones  $\alpha$  d'une chaîne polypeptidique dont la structure est lue dans la **PDB**.

Les chapitres 9 à 13 illustrent l'analyse et la comparaison de structures de protéines :

- diagramme de Ramachandran (chapitre 9) ;
- matrice des distances et carte des contacts d'une structure **3D** pour une protéine de la **PDB** (chapitre 10) ;
- superposition de structures protéiques et calcul du score **RMSD** à l'aide du module **Bio.PDB.Superimposer** de **Biopython** (chapitre 11) ;
- calcul de l'indice de repliement nommé *Foldindex* d'un ensemble de séquences de protéines lues dans un fichier **FASTA** à partir du site <https://fold.proteopedia.org> (chapitre 12) ;
- alignement de séquences et recherche de protéines homologues à l'aide de **PHI-BLAST** via le module **NCBIWWW** de **Biopython** (chapitre 13).

## Importance des questions en fin de chapitre

---

Chaque chapitre se termine par des questions dont les réponses sont fournies en fin d'ouvrage après le chapitre 13. L'ensemble des questions et réponses représente environ un cinquième de cet ouvrage : c'est un élément clé pour la réflexion du lecteur.

Ces questions et réponses constituent des rappels des notions développées dans chaque chapitre mais également des compléments et des illustrations. Elles permettent d'enrichir et de compléter les notions et les scripts tant du point de vue biologique qu'informatique. En effet, les scripts présentés dans le corps des chapitres sont courts et vont à l'essentiel.

Par exemple, au chapitre 1, pour obtenir un pourcentage, on effectue une division sans vérifier que le dénominateur n'est pas égal à zéro. Une des questions du chapitre 1

vient corriger ce problème et rappeler qu'il faut mettre systématiquement en place ce genre de vérifications.

De même, au chapitre 2, le script est censé lire une séquence dans un fichier dont le nom est passé en paramètre. Si le fichier n'existe pas, l'exécution du script aboutit à une erreur d'exécution. Ce comportement est volontaire et unique afin d'avoir un code `Python` court à lire. Une des questions du chapitre 2 a pour but de résoudre ce problème. Ensuite, tous les autres scripts qui utilisent un nom de fichier vérifient que le fichier existe en utilisant la technique indiquée dans la solution du chapitre 2.

## Publics visés

---

Cet ouvrage s'adresse à tous les étudiants, chercheurs, ingénieurs qui désirent approfondir leurs connaissances dans l'analyse de données biologiques avec des outils `Python` et `Biopython`. Les notions requises en biologie sont essentiellement celles relatives à la composition et à la structure tridimensionnelle des acides nucléiques et des protéines (nucléotides, acides aminés, chaînes polypeptidiques).

Les étudiants et étudiantes de premier cycle y trouveront une aide à l'utilisation de `Python` pour la bioinformatique tandis qu'il pourra servir de manuel de référence pour les étudiants de cycle supérieur, pour les chercheurs et ingénieurs qui débudent en bioinformatique et qui ont besoin de modèles de scripts pour gérer les séquences, interfacer les sites internationaux de bioinformatique, parcourir les fichiers de séquences, de structures `3D` ou les fichiers de résultats de requêtes, au format `JSON` ou `XML`.

Enfin, le dernier public concerné est celui, de plus en plus nombreux, de nos collègues enseignantes et enseignants en bioinformatique car ils trouveront au fil des chapitres des sujets de T.P. avec corrigés prêts à l'emploi et faciles à adapter ou compléter.

## Prérequis

---

La compréhension des scripts nécessite des connaissances élémentaires des instructions et fonctions du langage `Python`. Ces notions doivent également inclure les tableaux, les dictionnaires, les expressions régulières, les imports de packages et les appels de méthodes objets. Un rappel de ces connaissances `Python` se trouve dans les annexes, page 183, après le chapitre 13. Une présentation détaillée des expressions régulières et leur utilisation en `Python` est aussi fournie en annexe, page 191.

Aucune connaissance de `Biopython` n'est requise. Une présentation de `Biopython` est fournie en annexe (page 197). La liste des scripts écrits, avec leurs dépendances et la liste des modules et packages `Python` et `Biopython` utilisés sont également présentes dans cette annexe, respectivement pages 199 et 203.