

Rapport projet court

Sujet

Assignation des structures secondaire des protéines

Préparé par

Maha GRAA

Sommaire

1) Introduction :	3
2) Objectif et problématique :	3
3) Matériels et méthodes :	4
4) Résultats et discussion	6
5) Annexes	7
a. Lire fichier hb2 et faire le tri de la liste Hbond_final_tri_num (pièce centrale du script)	7
b. Fonction hélice	7
c. Fonction feuillet	8
d. Sortie en terminal	8

Projet : Assignation des structures secondaire des protéines

1) Introduction :

Au sein des cellules, on a trois classes majeurs de polymères qui sont :

- Acide ribonucléique (ARN)
- Acide désoxyribonucléique (ADN)
- Les protéines

C'est à la dernière classe qu'on va s'intéresser durant notre projet. En effet, une protéine est une chaîne d'acides aminés.

Les cellules du corps humain font appel à plus de vingt types d'acides aminés différents. On a des acides aminés essentiels pour le corps humain qui sont non synthétisables et qui doivent être apportés grâce à l'alimentation.

Ces acides aminés sont :

Valine, Leucine, Isoleucine, Lysine, Thréonine, Phénylalanine, méthionine, Histidine, tryptophane, Glutamine, Aspartate, Glutamate, Arginine, Alanine, Proline, Cystéine, Asparagine, Sérine, Glycine, Tyrosine. (Marc parisien. Éd (2005) Les feuillets beta dans les protéines).

Annotation, comparaison et construction :

Chaque acide aminé est composé d'atomes reliés entre eux par des liaisons covalentes. Ces atomes sont repartis en deux groupes : les atomes de la chaîne principale et les atomes de la chaîne latérale. Les atomes de la chaîne principale sont les mêmes pour tous les acides aminés.

Les acides aminés qu'on retrouve dans les protéines possèdent tous la même structure : qui est un atome de carbone central (α) sur lequel on trouve un groupe carboxyle (COOH) et un groupe aminé ($-\text{NH}_2$), un atome d'hydrogène ($-\text{H}$) et un groupement latéral (R). C'est grâce à ce groupement latéral que les acides aminés sont différents l'un de l'autre.

Il existe différents types de liaisons qui s'impliquent dans la structure des protéines : Structure primaire (liaisons peptidiques), secondaire (liaisons hydrogène) tertiaire et quaternaire.

On va s'intéresser en particulier à la structure secondaire, qui est obtenue grâce aux liaisons hydrogènes qui se trouvent au niveau des groupements $-\text{CO}$ et $-\text{NH}$ des liaisons peptidiques.

On doit se baser sur ces liaisons hydrogène pour pouvoir attribuer la forme de la structure secondaire d'une protéine. Ces structures peuvent être un hélices α (alpha) ou un feuillets β (bêta).

Il existe différentes méthodes pour qu'on puisse définir la structure secondaire citons comme exemple : STRIDE et DEFINE. L'algorithme DSSP (Define Secondary Structure of Proteins) reste la méthode standard pour attribuer une structure secondaire aux acides aminés.

DSSP est un algorithme ayant comme objectif le fait d'approcher la notion intuitive de structure secondaire par un algorithme objectif. Un algorithme qui sert à extraire des caractéristiques structurelles à partir des coordonnées atomiques d'où la reconnaissance des structures (Wiley & Sons (1983). Dictionary of Protein Secondary Structure: Pattern Recognition of Hydrogen-Bonded and Geometrical Features Biopolymers, Vol. 22, 2577-2637)

2) Objectif et problématique :

L'algorithme DSSP est un algorithme standard mais qui n'est pas toujours mis à jour (la dernière mise à jour a eu lieu en 2012).

Notre but durant ce projet est de trouver une méthode alternative à l'algorithme DSSP qui peut attribuer les structures secondaire à une protéine.

3) Matériels et méthodes :

Afin de pouvoir créer ce nouveau modèle :

On avait besoin d'installer Python sur notre machine de travail grâce à une installation de miniconda.

En effet Miniconda est une version minimaliste de Anaconda qui embarque Python ainsi que le gestionnaire de packages Conda et grâce à laquelle on peut ensuite installer des packages supplémentaires.

On avait besoin aussi d'un fichier PDB à partir de Protein Data Bank qui est La banque de données sur les protéines du Research Collaboratory for Structural Bioinformatics, plus communément appelée Protein Data Bank. Elle est une collection de données sur la structure tridimensionnelle de macromolécules biologiques qui sont les protéines essentiellement ainsi que les acides nucléiques. Les structures dans Data Bank sont déterminées par cristallographie aux rayons X ou grâce à la spectrométrie RMN. Les données obtenues suite à ces expériences seront déposées dans un fichier PDB. Le Format PDB est le format original de la banque. Ce format est pourtant relativement restrictif, il est en 80 colonnes. Le nombre maximum d'atomes d'un fichier pdb est de 99999. De même le nombre de résidus par chaîne est au maximum de 9999 : il n'y a que 4 colonnes autorisées pour ce chiffre. Le nombre de chaînes, lui, est limité à 62 : une seule colonne est disponible, et les valeurs possibles sont une des lettres des 26 lettres de l'alphabet, en minuscule ou en majuscule, ou un des chiffres de 0 à 9. Quand ce format a été déjà défini, ces limitations n'étaient pas restrictives, mais elles ont plusieurs fois été franchies lors du dépôt de structures extrêmement grandes citons les structures des virus, des ribosomes, des complexes multienzymatiques.

J'ai commencé par choisir une protéine d'intérêt, j'ai cherché dans Protein Data Bank sa structure et j'ai extrait le fichier contenant toutes les informations concernant cette protéine d'intérêt dans un fichier PDB. Pour mon projet, cette protéine est la tyrosine phosphatase yopH from yersinia pestis (codée 1huf dans PDB database). Les liaisons hydrogène ne sont pas citées dans le fichier PDB d'où l'obligation de passer par un programme qui permet de lire CES liaisons.

Parmi les solutions possibles, j'ai installé le programme HBPLUS (HBPLUS v.3.06 - a hydrogen bond calculation program).

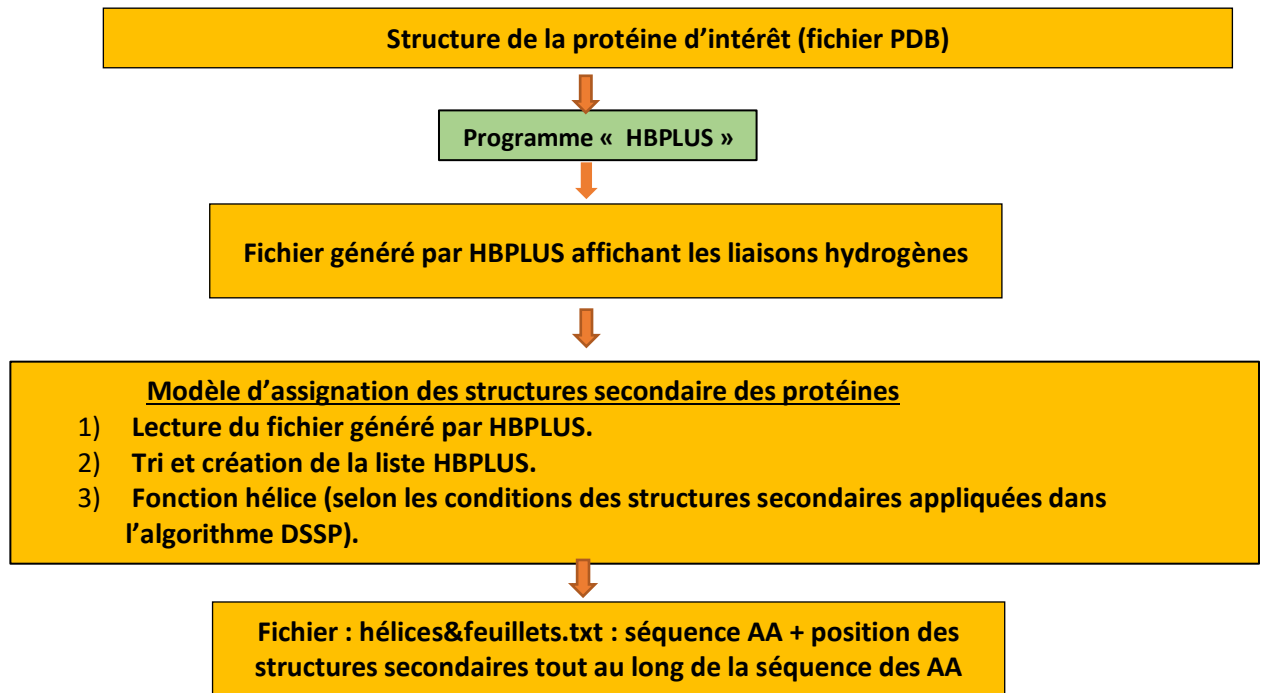
HBPLUS est un programme de calcul de liaisons hydrogène qui permet de calculer les géométries de toutes les liaisons hydrogène éventuellement des listes d'interactions avec les voisins, il calcule également les positions d'hydrogène, traite les hydrogènes qui peuvent occuper plus d'une position, les liaisons H amino-aromatiques, prend en charge la personnalisation complète, par exemple pour les liaisons H il décrit les types d'atomes donneurs et accepteurs.

Il Analyse la liaison H près des chaînes latérales Asn, Gln et His et suggère des conformations optimales

Son fichier de sortie est un fichier PDB comprenant les liaisons hydrogènes et leurs positions.

Ce fichier de sortie me permettra d'avancer dans la réalisation du projet. Une fois le fichier en question est récupéré, on pourra commencer à l'utiliser.

Ci-joint un workflow qui permet de décrire le fonctionnement du modèle :



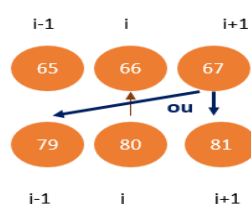
Description du Workflow (Modèle d'assignation des structures secondaire des protéines) :

Dans cette partie on va décrire le fonctionnement de la pièce centrale du workflow qui est le modèle principal :

- 1) Lecture du fichier généré par HBPLUS : pour pouvoir lire ce fichier, on a besoin d'un programme qui permet de lire ligne par ligne le fichier en question (matrice avec une seule colonne et plusieurs lignes) et attribuer les données en question à différentes listes qu'on utilisera ultérieurement.
- 2) Tri et création de la liste HBPLUS : créer une liste qui est dans ce cas la liste centrale du programme dans laquelle on va stocker toutes nos variables
- 3) Fonction hélice (selon les conditions des structures secondaire appliquées dans le programme DSSP) : au niveau de cette partie on a choisi d'appliquer les conditions suivies par le modèle standard DSSP, ces conditions sont bien détaillées dans un l'article (Wiley & Sons (1983). Dictionary of Protein Secondary Structure: Pattern Recognition of Hydrogen-Bonded and Geometrical Features Biopolymers, Vol. 22, 2577-2637)

Condition hélice : pour avoir une hélice, il faut avoir minimum deux turns de la même catégorie, les turns peuvent être des turns de 3, 4 ou 5. Cette condition est vraie si seulement si deux acides aminés qui se succèdent sont distants de 3, 4 ou 5 acides aminés. Une fois cette condition est validée, on vérifie si on a au moins deux turns de la même catégorie qui se succèdent. Si cette condition est validée on peut confirmer que dans cette partie de la structure on a une hélice. Il existe trois catégories d'hélices : 3 (hélice 3_{10}), 4 (hélice α), 5 (hélice π).

- 4) Fonction feuillet : la condition pour assigner la présence d'un feuillet bêta est la suivante : Si on a 3 acides aminés successifs qui ont Hbond avec autres 3 acides aminés, la condition qu'on supposera qu'elle est vraie est le fait que l'acide aminé au milieu (i) (voir figure) il a un Hbond avec un autre acide aminé (j) après on passe à la vérification de l'acide aminé suivant (i+1) s'il a un Hbond avec l'acide aminé qui vient juste après (j+1) ou juste avant (j-1) celui qui a un Hbond avec notre condition vraie du début. Ci-dessous une figure qui explique les conditions appliquées.



4) Résultats et discussion

En utilisant le modèle mis en place, on aura comme résultat un fichier txt intitulé helices&feuillets qui contient 6 colonnes ([a]position de l'acide aminé, [Type]le nom de l'acide aminé, [hélice]position des hélices, [turn]position et type 3,4, ou 5 des Turn, [bridge]position des bridges,[leader]position des leaders et [feuille] position des feuillets.

Si on prend comme exemple une protéine modèle « TYROSINE PHOSPHATASE YOPH FROM YERSINIA PESTIS », ayant comme code 1HUF dans la base protein PDB, selon notre modèle, on trouvera 6 hélices et 4 feuillets.

Faisons une matrice de confusion pour comparer le modèle standard DSSP avec notre modèle (sans séparer hélices et feuillets dans la matrice) :

True positive (7)	False negative (3)
False positive (4)	True negative ?

On peut constater que sur 10 structures secondaires citées dans le DSSP, on a pu identifier 7 structures secondaire grâce à notre modèle. Notre précision est à 73,33 %

Ainsi on peut constater que parmi 10 structures secondaires retrouvées grâce notre modèle, 7 structures sont vrais.


```
helice['-', '-', '-', '-', '-', '-', '-', 'H', 'H', '-', '-', '-', '-', 'H', 'H', 'H', 'H', 'H', 'H', 'H', '-', 'H', 'H
```

c. Fonction feuillet

```

61  #*****Fonction feuillet *****
62  i=1
63  for i in range (longueur-1) :
64      a = int (Hbond_final_trie_num[i][0])
65      b = int (Hbond_final_trie_num[i][1])
66      a1 = int (Hbond_final_trie_num[i-1][0])
67      a2 = int (Hbond_final_trie_num[i+1][0])
68      b2 = int (Hbond_final_trie_num[i+1][1])
69      if a1<a and a<a2 and a2-a == 1 and a - a1 == 1 :
70          if abs (b-b2) == 1 :          # utiliser le b de i+1
71              bridges [i] = "B"
72
73  i=0
74  for i in range (longueur):
75      if bridges[i] == "B" and bridges[i-1]=="B" :
76          leader[i-1] = "L"
77          i=i+1
78  i=0
79  for i in range (longueur):
80      if leader[i] == "L" and leader[i-1]=="L" and feuillet[i-1] == '-' :
81          feuillet [i] = "F"
82          feuillet[i-1] = "F"
83      i=i+2

```

[illegible]

d. Sortie en terminal

```

90 #***** Sortie en Terminal *****
91
92 print("      a      type  helice  turn   bridge leader feuillet")
93 for i in range (longueur) :
94     print("      ", Hbond_final_trie_num[i][0], "      ", type[i] , "      ", helices[i], "      ", turn [i], "      ", bridges[i] , "      ", leader[i] , "      ", feuillet[i] )
95
96
97 with open("helices&feuillets.txt", 'w') as f:
98     f.write("      a      type  helice  turn   bridge leader feuillet\n")
99     for i in range (longueur) :

```

a	type	helice	turn	bridge	leader	feuillet
0003	LEU	-	-	-	-	-
0004	SER	-	-	-	-	-
0004	SER	-	-	-	-	-
0005	LEU	-	-	-	-	-
0006	SER	-	-	-	-	-
0006	SER	-	-	-	-	-
0007	ASP	-	3	-	-	-
0008	LEU	H	4	B	-	-
0009	HIS	H	4	-	-	-
0009	HIS	-	-	-	-	-
0009	HIS	-	-	-	-	-
0010	ARG	-	4	-	-	-
0010	ARG	-	3	-	-	-
0011	GLN	H	4	B	L	-
0012	VAL	H	4	B	-	-
0013	SER	H	4	-	-	-
0013	SER	H	4	-	-	-
0013	SER	H	4	-	-	-
0014	ARG	H	4	-	-	-
0014	ARG	H	4	-	-	-
0014	ARG	-	-	-	-	-
0015	LEU	H	4	B	L	F
0016	VAL	H	4	B	L	F
0017	GLN	H	4	B	-	-
0018	GLN	H	4	-	-	-
0018	GLN	-	-	-	-	-
0019	GLU	-	3	B	-	-
0020	SER	-	5	-	-	-
0020	SER	-	-	-	-	-