

INFO-F-208 - Projet 1 -- Alexandre Heneffe -- 12/10/2018

Ce document "jupyter notebook" nous introduit et nous décrit les différents algorithmes implémentés dans le cadre du Mini Projet 1 du cours "Introduction à la bioinformatique"

Introduction

Pour ce premier Mini Projet nous traitons la problématique liée à l'alignement de protéines constituées de séquence d'acides aminés représentées par des lettres. Les connaissances nécessaires à la compréhension de ce projet sont les bases de la biologie (Notion de protéine, acides aminés).

Nous procéderons à l'alignement de séquences de protéines mais également de BRD ("Bromo Domain"). Les bromodomains sont des modules interactifs reconnaissant les sites d'acétylation (réactions chimiques) dans les protéines. Les BRD sont notamment intéressantes pour développer de nouvelles thérapies pour le cancer humain ou pour la sclérose en plaques par exemple.

Les acides aminés ont évolué au fil du temps, ce qui a induit des modifications (insertions, suppressions, substitutions). De ce fait, nous cherchons à les aligner de sorte à trouver leurs similitudes maximales dissimulées et déterminer si celles-ci sont homologues (2 gènes ou 2 protéines qui partagent un ancêtre commun).

Notre objectif est d'implémenter les algorithmes de Needleman-Wunsch et de Smith-Waterman qui trouveront respectivement les alignements globaux et les alignements locaux et permettront d'analyser les similitudes. Ces notions seront détaillées plus tard.

Matériel

Nous allons maintenant aborder l'implémentation des ADT "Abstract Data Type" permettant de créer des outils utiles à l'implémentation des algorithmes.

```
In [1]: 1 import copy as c
        2 import pandas
        3
        4 pandas.set_option('display.max_columns', 500)
```

ADT Séquence

Avec cette classe, nous créons des objets représentant des Séquences. Nous pouvons l'afficher au format FASTA pour visualiser son nom et ses acides aminées par exemple.

```

In [2]: 1 class Sequence:
2         """
3         Classe qui représente un objet Séquence d'acides aminées
4         """
5
6         def __init__(self, title = None, seq = None):
7             """ Crée un objet Séquence """
8
9             if title != None and seq != None:
10                self.__seq = seq # Séquence : string
11                self.__title = title #Titre de la séquence
12
13         def get_acids(self):
14             """ Renvoie les lettres de la séquence """
15             return self.__seq
16
17         def length(self):
18             """ Renvoie la taille de la séquence """
19             return len(self.__seq)
20
21         def set_title(self, title):
22             """ Donne un nom a la séquence """
23             self.__title = title
24
25         def set_seq(self, seq):
26             """ met à jour la séquence """
27             self.__seq = seq
28
29         def display(self):
30             """ Affiche la séquence au format FASTA """
31             print(self.__title)
32             print(self.__seq, end="\n\n")
33
34
35         print("Exemple de séquence: ", end= 2*"\\n")
36         ExempleSeq = Sequence("sp|060885|75-147", "WKHQFAWPFQPPVDAVKLNLPDYKIIKTPMDMGTIKKRLENNYYWNAQECIQDFNTMFTNCYIYNKPGDDIV")
37         ExempleSeq.display()
38
39
40

```

Exemple de séquence:

```

sp|060885|75-147
WKHQFAWPFQPPVDAVKLNLPDYKIIKTPMDMGTIKKRLENNYYWNAQECIQDFNTMFTNCYIYNKPGDDIV

```

ADT Parser de séquence(s)

Nous utilisons le parser de séquence(s) afin de récupérer les séquences d'acides aminées ainsi que leurs noms depuis un fichier au format FASTA (succession de séquences précédées par leurs noms). Nous créons ensuite une liste d'objets Séquence à partir des données fournies que nous pourrions utiliser dans les algorithmes décrits plus tard.

```

In [3]: 1 class ParserSequence:
2         """
3         Classe qui représente un Parser qui va lire un fichier avec des séquences
4         """
5
6         def __init__(self, file_name):
7             """ Crée un objet Parser """
8
9             self.lines = open(file_name, "r").readlines() #Lignes du fichier
10            self.nb_lines = len(self.lines)
11            self.sequences = [] #Liste d'objets Séquence
12
13            def create_seq(self, seq, title):
14                """ Ajoute une séquence a la liste de séquences """
15
16                new_S = Sequence() # Crée un nouvel objet Séquence
17                new_S.set_seq(seq) # lui donne la séquence correspondante
18                new_S.set_title(title) # titre de la séquence
19
20                self.sequences.append(new_S) #ajout d'un nouvel objet séquence
21
22            def get_seq(self, i):
23                """ Renvoie une séquence de la liste """
24                return self.sequences[i]
25
26            def parse(self):
27                """ Récupère les séquences du fichier """
28
29                tmp = ""
30                title = ""
31                for line in self.lines: # lignes du fichier
32                    if line[0] == ">":
33                        # print(tmp + "\n")
34                        if tmp != "":
35                            self.create_seq(tmp, title) # on crée un objet séquence
36                            tmp = ""
37
38                            title = line
39                        else:
40                            tmp += line.strip("\n")
41
42                self.create_seq(tmp, title)

```

ADT Matrice

Afin d'aligner les séquences d'acides aminées, nous utilisons des objets représentant des matrices. Les différents types de matrice qui dérivent de cet objet Matrice seront décrits après (matrice de substitution, de scoring, V, W).

```

In [4]: 1 class Matrix:
2         """
3         Classe qui représente un objet Matrice
4         """
5
6         def __init__(self):
7             """ Crée un objet Matrice """
8
9             self.mat = []
10            self.n = 0 # colonnes (seq1)
11            self.m = 0 # lignes (seq2)
12            self.letters_seq1 = {} # dictionnaire clé = lettre, valeur = index
13            self.letters_seq2 = {} # "
14            self.letters_orders1 = "" #Les lettres des colonnes mis dans l'ordre
15            self.letters_orders2 = "" #Les lettres des lignes mis dans l'ordre
16
17
18            def get_acid_score(self, i, j):
19                """ Renvoie le score d'une cellule représentée par 2 lettres """
20                if i == "-" or j == "-":
21                    return -1
22                else:
23                    return self.mat[self.letters_seq1[i]][self.letters_seq2[j]]
24
25            def get_score(self, i, j):
26                """ Renvoie le score d'une cellule en mat[i][j] """
27                return self.mat[i][j]
28
29            def set_score(self, i, j, score):
30                """ donne un score a une cellule """
31                self.mat[i][j] = score
32
33            def addline(self):
34                """ ajoute une ligne dans la matrice """
35                self.mat.append([])
36                self.m += 1
37
38            def add_cell(self, i, score):
39                """ ajoute un élément dans une ligne (score) """
40                if isinstance(score, str):
41                    score = int(score)
42
43                self.mat[i].append(score)
44
45            def set_lign(self, scores):
46                """ ajoute toute une ligne avec des scores """
47                self.mat.append(scores)
48
49            def inc_nb_col(self):
50                """ augmente le nb de colonnes de 1 """
51                self.n += 1
52
53            def set_nb_col(self):
54                """ détermine le nombre de colonnes """
55                self.n = len(self.mat[0])
56
57            def inc_nb_lign(self):
58                """ augmente le nb de lignes de 1 """
59                self.m += 1
60

```

La matrice est déterminée par 2 séquences de lettres. Nous définissons une structure de la forme:

$$\{ " A ": 0, " B ": 1, " C ": 2, \dots \}$$

Cela nous permettra de récupérer le score d'une paire d'acides aminés à partir de leurs lettres. Dans notre cas, ce sera pour récupérer des scores figurant dans la matrice de substitution.

```
In [5]: 1 def get_letters(self, i, j):
2         """ Renvoie les lettres correspondantes a la cellule i,j """
3         return self.letters_orders2[i] + self.letters_orders1[j]
4
5     def set_letters_seq(self, seq1, seq2):
6         """ remplit les dictionnaires par
7             clé : acide aminée, valeur : indice dans la matrice
8         """
9
10        seq1 = seq1.replace(" ", "")
11        seq2 = seq2.replace(" ", "")
12        self.letters_orders1 = seq1
13        self.letters_orders2 = seq2
14
15        self.letters_seq1 = dict(zip(seq1, (i for i in range(len(seq1)))))
16        self.letters_seq2 = dict(zip(seq2, (i for i in range(len(seq2)))))
17
18        Matrix.get_letters = get_letters
19        Matrix.set_letters_seq = set_letters_seq
```

Pour visualiser les matrices tout au long des algorithmes, nous les affichons de manière jolie à l'aide de la méthode ci-dessous.

```
In [6]: 1 def display(self):
2         """ Affiche la matrice """
3
4         for l in self.letters_orders1:
5             if l == "-":
6                 print("\t-", end="\t")
7             else:
8                 print(l, end=" \t")
9         print("\n")
10
11        for i in range(self.m): # lignes
12            print(self.letters_orders2[i], end="\t")
13            if self.letters_orders2[i] == "-":
14                #print(" ", end=" ")
15                pass
16            for j in range(self.n): # colonne
17                print(self.mat[i][j], end=" \t")
18            print("\n")
19        print("\n")
20
21        Matrix.display = display
22
23    def panda(self):
24        """ Autre moyen d'afficher la matrice de manière plus compacte """
25
26        print(pandas.DataFrame(self.mat, list(self.letters_orders2), \
27                                pandas.Index(list(self.letters_orders1), name="*")), end="\n\n")
28
29        Matrix.panda = panda
30
```

ADT Matrice de substitution

La matrice de substitution nous donne le score de similarité de chaque paire d'acides aminés.

Lorsque 2 acides aminés se substituent, l'évolution accepte ou non le remplacement selon la préservation de la stabilité de la protéine et la modification de leurs fonctionnalités biochimiques. Le score de similarité nous donne une idée de ces remplacements en nous donnant la probabilité que les 2 acides aminés soient homologues et la probabilité qu'ils soient liés de façon stochastique. Un score positif montre les mutations acceptées et un score négatif montre les mutations refusées.

Nous avons à notre disposition les matrices de substitution "BLOSUM62", "BLOSUM80", "PAM120 et "PAM60" dont les scores ont été calculés sur base de l'homologie de différentes espèces. Les différents types de matrices "PAM" montrent une différence dans le nombre de substitutions acceptées, donc la divergence entre les séquences tandis que les différents types de "BLOSUM" sont dues à une variation du pourcentage d'identité entre les séquences.

```
In [7]: 1 class MatSubstitution(Matrix):
2         """
3         Classe qui représente un Parser qui va lire un fichier avec une matrice
4         de substitution
5         Cette classe hérite de Matrix
6         """
7
8         def __init__(self, file_name):
9             """ Crée un objet Parser """
10
11             Matrix.__init__(self)
12             self.lines = open(file_name, "r").readlines() #Lignes du fichier
13             self.nb_lines = len(self.lines)
14
15         def get_mat_sub(self):
16             """ Renvoie la matrice de substitution parsée """
17
18             return self.mat
19
```

ADT Parser de matrice

Nous utilisons le parser de matrice pour récupérer une matrice de substitution contenue dans un fichier et la placer dans un objet "Matrice" sur lequel nous pouvons appliquer différentes fonctions (tel que l'afficher, obtenir l'information d'une cellule, ajouter des éléments, etc..) . Nous nous servirons des scores de cette matrice de substitution lors de la détermination des alignements de séquences d'acides aminés.

```

In [8]: 1 def parse(self):
        2     """ Récupère la matrice du fichier """
        3
        4     i = 0
        5     for line in self.lines: # lignes du fichier
        6
        7         line = line.strip("\n")
        8         if len(line) > 0:
        9             if line[0] == " ":
        10                 self.set_letters_seq(line, line) # lettres de seq 1 et
        11             else:
        12                 if line[0] != "#":
        13                     self.addline()
        14                     all_line = line.split(" ")
        15                     for l in all_line:
        16                         if l != "" and not l.isalpha() and l != "*":
        17                             self.add_cell(i, l)
        18                     i += 1
        19
        20         self.set_nb_col()
        21
        22 MatSubstitution.parse = parse
        23
        24 print("Matrice de substitution BLOSUM62: \n\n")
        25 subEx = MatSubstitution("blosum62.txt")
        26 subEx.parse()
        27 #subEx.display()
        28 subEx.panda()
        29

```

Matrice de substitution BLOSUM62:

```

*  A  R  N  D  C  Q  E  G  H  I  L  K  M  F  P  S  T  W  Y  V  B  Z  X
*
A  4 -1 -2 -2  0 -1 -1  0 -2 -1 -1 -1 -1 -2 -1  1  0 -3 -2  0 -2 -1  0 -
4
R -1  5  0 -2 -3  1  0 -2  0 -3 -2  2 -1 -3 -2 -1 -1 -3 -2 -3 -1  0 -1 -
4
N -2  0  6  1 -3  0  0  0  1 -3 -3  0 -2 -3 -2  1  0 -4 -2 -3  3  0 -1 -
4
D -2 -2  1  6 -3  0  2 -1 -1 -3 -4 -1 -3 -3 -1  0 -1 -4 -3 -3  4  1 -1 -
4
C  0 -3 -3 -3  9 -3 -4 -3 -3 -1 -1 -3 -1 -2 -3 -1 -1 -2 -2 -1 -3 -3 -2 -
4
Q -1  1  0  0 -3  5  2 -2  0 -3 -2  1  0 -3 -1  0 -1 -2 -1 -2  0  3 -1 -
4
E -1  0  0  2 -4  2  5 -2  0 -3 -3  1 -2 -3 -1  0 -1 -3 -2 -2  1  4 -1 -
4
G  0 -2  0 -1 -3 -2 -2  6 -2 -4 -4 -2 -3 -3 -2  0 -2 -2 -3 -3 -1 -2 -1 -
4
H -2  0  1 -1 -3  0  0 -2  8 -3 -3 -1 -2 -1 -2 -1 -2 -2  2 -3  0  0 -1 -
4
I -1 -3 -3 -3 -1 -3 -3 -4 -3  4  2 -3  1  0 -3 -2 -1 -3 -1  3 -3 -3 -1 -
4
L -1 -2 -3 -4 -1 -2 -3 -4 -3  2  4 -2  2  0 -3 -2 -1 -2 -1  1 -4 -3 -1 -
4
K -1  2  0 -1 -3  1  1 -2 -1 -3 -2  5 -1 -3 -1  0 -1 -3 -2 -2  0  1 -1 -
4
M -1 -1 -2 -3 -1  0 -2 -3 -2  1  2 -1  5  0 -2 -1 -1 -1 -1  1 -3 -1 -1 -
4
F -2 -3 -3 -3 -2 -3 -3 -3 -1  0  0 -3  0  6 -4 -2 -2  1  3 -1 -3 -3 -1 -
4
P -1 -2 -2 -1 -3 -1 -1 -2 -2 -3 -3 -1 -2 -4  7 -1 -1 -4 -3 -2 -2 -1 -2 -
4
S  1 -1  1  0 -1  0  0  0 -1 -2 -2  0 -1 -2 -1  4  1 -3 -2 -2  0  0  0 -
4
T  0 -1  0 -1 -1 -1 -1 -2 -2 -1 -1 -1 -1 -2 -1  1  5 -2 -2  0 -1 -1  0 -
4

```

ADT Matrice de Scoring

Nous avons une classe qui permet de créer une matrice de scoring et de l'initialiser selon la méthode. Si nous sommes dans un alignement global, nous initialisons la matrice avec les pénalités d'ouverture de gap "I" et de continuité de gap "E". La notion de "gap" sera détaillée plus tard.

```
In [9]: 1 class MatScoring(Matrix):
2         """ Classe qui représente une matrice contenant les scores avec une p
3         de gap affine
4         Cette classe hérite de Matrix
5         """
6
7         def __init__(self, I, E, n, m):
8             Matrix.__init__(self)
9
10            self.n = n
11            self.m = m
12            self.I = I
13            self.E = E
14
15            def init_S_global(self):
16                """ Initialise lere ligne et lere colonne de S pour la méthode gl
17
18                first_col = [0]
19                for i in range(self.m):
20
21                    newline = []
22                    if i == 0:
23                        newline += [0, -self.I] # lere ligne [0, -I, -I-E, ...
24                        for j in range(2, self.n):
25                            newline.append(newline[j-1] - self.E)
26
27
28                    elif i == 1:
29                        newline += [-self.I] + (self.n-1)*[""]
30                        first_col += [-self.I]
31                    else:
32                        first_col.append(first_col[i-1] - self.E)
33                        newline += [first_col[i]] + (self.n-1)*[""]
34
35                    self.set_lign(newline)
36
37            # ===== EXEMPLE =====
38            ExempleS = MatScoring(4, 2, 5, 4)
39            ExempleS.init_S_global()
40            ExempleS.set_letters_seq("-ABCD", "-EFGH")
41            print("Exemple de matrice de scoring initialisée pour l'alignement global")
42            ExempleS.display()
```

Exemple de matrice de scoring initialisée pour l'alignement global:

	-	A	B	C	D
-	0	-4	-6	-8	-10
E	-4				
F	-6				
G	-8				

ADT Matrices V et W

Nous pouvons créer des matrices V et W permettant de sauvegarder les informations concernant les gap lors du calcul de la matrice de scoring. Nous les initialisons avec des valeurs infinies - Inf sur la première ligne de V et des 0 sur la première colonne de V. Nous procédons de la même manière pour W mais en inversant les valeurs -Inf et 0.

$$V = \begin{bmatrix} -Inf & -Inf & \dots \\ 0 & & \\ 0 & & \\ \dots & & \end{bmatrix} \quad W = \begin{bmatrix} 0 & 0 & \dots \\ -Inf & & \\ -Inf & & \\ \dots & & \end{bmatrix}$$

```
In [10]: 1 class MatV(Matrix):
2         """ Classe qui représente une matrice permettant de sauvegarder des v
3           liées aux lignes lors du calcul de la matrice de scoring
4           Cette classe hérite de Matrix
5         """
6
7         def __init__(self, n, m):
8             Matrix.__init__(self)
9
10            self.n = n
11            self.m = m
12
13            def init_V(self):
14                """ Initialise lere ligne et lere colonne de V """
15
16                for i in range(self.m):
17                    if i == 0:
18                        self.set_lign([-float("inf")]*self.n) # -inf -inf -inf ..
19                    else:
20                        self.set_lign([0] + (self.n-1)*[""]) # 0 ...
```

```
In [11]: 1 class MatW(Matrix):
2         """ Classe qui représente une matrice permettant de sauvegarder des v
3           liées aux colonnes lors du calcul de la matrice de scoring
4           Cette classe hérite de Matrix
5         """
6
7         def __init__(self, n, m):
8             Matrix.__init__(self)
9
10            self.n = n
11            self.m = m
12
13            def init_W(self):
14                """ Initialise lere ligne et lere colonne de W """
15
16                for i in range(self.m):
17                    if i == 0:
18                        self.set_lign([-float("inf")] + [0]*(self.n-1)) # -inf 0
19                    else:
20                        self.set_lign([-float("inf")] + (self.n-1)*[""]) # 0 ...
```

Méthodes

Maintenant, nous allons parler de la méthodologie utilisée pour produire des résultats corrects.

Objectif

Ce que nous cherchons à obtenir à l'aide de la programmation dynamique et les algorithmes, ce sont des solutions de cette forme :

```
In [12]: 1 print("MGGETFA\n :: : \n-GGVTTTF")
```

```
MGGETFA
:: :
-GGVTTTF
```

Ceci est un exemple d'alignement de 2 séquences d'acides aminées.

Nous introduisons la notion de "gap". Celui-ci est un espace dans la séquence qui est représenté par un "-". Il sert à procéder à une insertion ou à une suppression. Dans l'exemple ci-dessus, l'élément "M" de la première séquence est supprimé dans la deuxième séquence. Cela permet de ne garder que les similitudes entre les 2 séquences.

De plus, étant donné que nous utilisons un système de scores pour procéder à ces alignements de séquences, et que les gaps diminuent les pourcentages de similarité entre 2 séquences, nous leur donnons un score de pénalité.

Dans le cadre de ce projet, nous utiliserons une pénalité de gap affine et non linéaire:

- Une pénalité linéaire est une valeur négative identique pour tous les gaps
- Avec une pénalité affine, une ouverture de gap est plus pénalisante qu'une continuité de gap:

$$ABC - - - -DF$$

Par exemple, le premier gap aura une pénalité de -4 (I) et ceux qui suivent auront une pénalité de -1 (E).

Programmation dynamique

Le meilleur moyen de trouver les alignements optimaux est la programmation dynamique. En effet, dans une approche de "Diviser pour régner", un résultat à l'instant t est déterminé par un résultat à l'instant $t-1$. Si nous devons calculer toutes les combinaisons possibles entre les 2 séquences d'aminées, le temps de calcul serait énorme. A la place, à chaque étape, nous utilisons les résultats précédemment calculés.

Notre problème possède les 2 caractéristiques représentatives d'un problème devant être résolu par de la programmation dynamique:

- Le problème contient des solutions optimales pour les sous-problèmes du problème global: nous retenons les valeurs $V(i,j)$ et $W(i,j)$ à chaque étape et $S(i,j)$ détermine un sous-alignement de séquences optimales $seq1[j:]$ et $seq[i:]$ (des sous-séquences).
- La solution récursive contient un nombre limité de sous-problèmes distincts répétés beaucoup de fois: Lors du backtracking, nous construisons d'abord la solution optimale de taille minimale (en bas à droite) et puis nous construisons de manière croissante les solutions optimales des sous-problèmes en remontant vers le coin en haut à gauche.

Notre cas ressemble fortement à un problème de recherche de distance d'édition "D" entre 2 chaînes de caractères (Nombre minimum d'opérations d'insertion, suppression, remplacement pour passer d'une chaîne à une autre).

En effet, 2 caractères différents dans un alignement de séquence représentent un remplacement, un espace dans la première séquence représente une insertion dans la seconde et un espace dans la seconde séquence représente une suppression dans la première séquence.

il possède également les 2 caractéristiques d'un problème dynamique et réalise une solution récursive de la forme:

$$D(i,j) = \begin{cases} i & \text{si } \{i > 0, j = 0\} \\ j & \text{si } \{i = 0, j > 0\} \\ \min(D(i-1,j) + 1, D(i,j-1) + 1, D(i-1,j-1) + t(i,j)) & \text{si } \{i > 0, j > 0\} \end{cases}$$

similaire à celle de l'alignement de séquence (détaillé plus tard)

Alignement

Nous définissons un objet Alignement qui représente un alignement de séquence de manière générale. Celui-ci se décline en 2 types d'Alignements:

- L'alignement global
- L'alignement local

Tous les 2 ont des algorithmes différents mais similaires d'un point de vue technique. L'alignement général va simplement initialiser les séquences de séquences d'acides aminés, les matrices de substitution, de scoring, les matrices de sauvegarde V et W pour les alignements globaux ou locaux.

```

In [13]: 1 class Alignment:
2         """ Classe représentant un objet qui trouvera l'alignement de 2 séquences
3           d'acides aminées
4         """
5
6         def __init__(self, I, E, mat_file, seq1, seq2, p):
7
8             self.I = I
9             self.E = E
10            self.p = p
11            self.seq1 = seq1
12            self.seq2 = seq2
13            self.n = self.seq1.length()+1
14            self.m = self.seq2.length()+1
15
16            if self.p == 1 or p == 2: # Si on veut imprimer les informations
17                print("Séquence 1 de longueur {0}: ".format(self.n))
18                self.seq1.display()
19                print("Séquence 2 de longueur {0}: ".format(self.m))
20                self.seq2.display()
21                print("matrice de substitution utilisée: {0}".format(mat_file))
22                print("Pénalité de gap affine: I = {0} | E = {1}".format(self.I, self.E))
23
24            # ===== SUBSTITUTION =====
25            self.t = MatSubstitution(mat_file)
26            self.t.parse()
27
28            # ===== SCORING =====
29
30            self.S = MatScoring(I, E, self.n, self.m)
31            # on crée un objet matrice de scoring
32
33            # ===== V ET W =====
34
35            self.V = MatV(self.n, self.m)
36            self.W = MatW(self.n, self.m)
37
38            self.V.init_V() # Initialise V
39            self.V.set_letters_seq("-"+self.seq1.get_acids(), "-"+self.seq2.get_acids())
40            self.W.init_W() # Initialise W
41            self.W.set_letters_seq("-"+self.seq1.get_acids(), "-"+self.seq2.get_acids())
42
43            self.current_sol = [] # Pour le backtracking
44            self.all_solutions = []
45
46            def get_v(self, i, j):
47                return max( self.S.get_score(i-1, j) - self.I,
48                           self.V.get_score(i-1,j) - self.E )
49
50            def get_w(self, i, j):
51                return max( self.S.get_score(i, j-1) - self.I,
52                           self.W.get_score(i, j-1) - self.E )

```

Alignement global

L'alignement global va utiliser l'algorithme de Needleman et Wunsch (implémentée ci-dessous dans la classe Alignement Global). Celui-ci va chercher l'alignement qui donne la plus grande similarité entre 2 séquences d'acides aminées.

Tout d'abord, Nous devons initialiser la matrice de scoring avec une pénalité de gap affine. Comme décrit plus haut, la matrice initiale aura la forme:

$$S = \begin{bmatrix} 0 & -I & -I - E & -I - E - E & \dots \\ -I & & & & \\ -I - E & & & & \\ -I - E - E & & & & \\ \dots & & & & \end{bmatrix}$$

Ensuite, nous calculons le reste des cellules de cette matrice:

Etant donné que des gaps sont introduits dans les 2 séquences, il faut sauvegarder l'état de ces gaps dans des matrices V et W préalablement initialisée comme décrit plus haut.

- Première séquence (les colonnes):
 - Si la valeur précédente N'était PAS un gap, on utilise $S(i - 1, j)$ - la pénalité de gap d'ouverture I
 - Si la valeur précédente ETAIT un gap, on utilise $V(i - 1, j)$ - la pénalité de gap de continuité E

$$V(i, j) = \max \begin{cases} S(i - 1, j) - I \\ V(i - 1, j) - E \end{cases}$$

- Deuxième séquence (les lignes):
 - Si la valeur précédente N'était PAS un gap, on utilise $S(i, j - 1)$ - la pénalité de gap d'ouverture I
 - Si la valeur précédente ETAIT un gap, on utilise $W(i, j - 1)$ - la pénalité de gap de continuité E

$$W(i, j) = \max \begin{cases} S(i, j - 1) - I \\ W(i, j - 1) - E \end{cases}$$

La cellule $S(i, j)$ de la matrice de scoring devient alors:

$$S(i, j) = \max \begin{cases} S(i - 1, j - 1) + MatSubstitution(i, j) \\ V(i, j) \\ W(i, j) \end{cases}$$

```

In [14]: 1 class GlobalAlignment(Alignment):
2         """
3         Classe qui va trouver un aligement de séquences d'acides aminées avec
4         une pénalité affine et une méthode globale
5         """
6
7     def __init__(self, k, I, E, mat_file, seq1, seq2, p):
8         Alignment.__init__(self, I, E, mat_file, seq1, seq2, p)
9         self.k = k
10
11         self.S.init_S_global() # On initialise la matrice de scoring
12         self.S.set_letters_seq("-"+self.seq1.get_acids(), "-"+self.seq2.get_acids())
13
14     def update_s_global(self, i, j, v_ij, w_ij):
15         letters_ij = self.S.get_letters(i,j) # 'AB' par exemple
16         t_ij = self.t.get_acid_score(letters_ij[0], letters_ij[1])
17
18         s_ij = max( self.S.get_score(i-1,j-1) + t_ij, v_ij, w_ij )
19         self.S.set_score(i,j, s_ij)
20
21     def Needleman_Wunsch(self):
22         """
23         Algorithme qui calcule la matrice de scoring pour l'alignement global
24         en utilisant la pénalité affine puis fait un backtracking pour récupérer
25         tous les alignements optimaux possibles
26         """
27
28         # ===== CREATION SCORING =====
29
30         for i in range(1, self.m):
31             for j in range(1, self.n):
32                 v_ij = self.get_v(i, j) # V(i,j)
33                 self.V.set_score(i,j, v_ij)
34                 w_ij = self.get_w(i, j)
35                 self.W.set_score(i,j, w_ij) # W(i,j)
36
37                 self.update_s_global(i, j, v_ij, w_ij)
38
39         if self.p == 2: # Si on veut afficher les matrices
40             self.S.panda()
41
42         self.current_sol.append((self.m-1, self.n-1))
43         self.backtracking_global(self.m-1, self.n-1) #appel sur element de la matrice
44                                                         #ligne, derniere colonne
45
46         R = Result(self.S, self.t, self.all_solutions, self.p)
47         res = R.compute_result()
48         return res
49

```

backtracking

Une fois la matrice de scoring créée, il faut utiliser un backtracking. Le backtracking est une méthode algorithmique permettant de trouver toutes les solutions optimales d'un problème de manière récursive. Il construit une solution petit à petit selon certaines conditions. Une fois qu'il a trouvé une solution admissible, il détruit la solution et en crée une nouvelle. C'est une approche "BOTTOM-UP".

Dans notre cas, le backtracking part du dernier élément de la matrice de scoring (en bas à droite) et va remonter vers l'élément en (0,0). Une solution est donc un chemin partant du dernier élément jusqu'au premier élément.

A chaque étape nous avons 3 possibilités:

$$\begin{bmatrix} S(i-1, j-1) + t(i, j) & V(i, j) \\ & W(i, j) & S(i, j) \end{bmatrix}$$

Etant donné que S(i,j) est calculée à partir de V, W, et de l'élément en diagonale + le score de la matrice de substitution, il revient sur la cellule ayant le score maximal de ces 3 éléments.

La méthode ci-dessous réalise ce backtracking:

```
In [15]: 1 def backtracking_global(self, i, j):
2         """ Remonte la matrice de scoring a partir du dernier élément jusqu'à
3             pour avoir les k alignements
4             """
5
6         if i == 0 or j == 0:
7             if i == 0 and j == 0:
8                 if len(self.current_sol) > 0:
9                     self.current_sol.pop() # Si on est en (0,0)
10                if len(self.all_solutions) == self.k: # Si on a déjà trouvé k ali
11                    return
12                if self.current_sol not in self.all_solutions:
13                    #print("1 solution trouvé: ", self.current_sol)
14                    self.all_solutions.append(c.deepcopy(self.current_sol))
15
16            else:
17                for pos in range(3):
18
19                    new_i = i
20                    new_j = j
21                    valid = False
22                    if pos == 0 and self.is_previous("v", i, j): # haut
23                        new_i -= 1 # i-1
24                        valid = True
25                    elif pos == 1 and self.is_previous("w", i, j): # gauche
26                        new_j -= 1 # j-1
27                        valid = True
28                    elif pos == 2 and self.is_previous("s", i, j): # diagonale
29                        new_i -= 1 # i - 1
30                        new_j -= 1 # j - 1
31                        valid = True
32
33                    if valid:
34                        self.current_sol.append((new_i, new_j))
35                        self.backtracking_global(new_i, new_j) # appel sur cellule
36                        if len(self.current_sol) != 0:
37                            self.current_sol.pop() # destruction sol partielle
38
39 GlobalAlignment.backtracking_global = backtracking_global
```

La méthode `is_previous` (ci-dessous) permet à tout moment lors du backtracking de vérifier vers quelle cellule il doit se diriger:

- Si l'élément $S(i,j)$ est déterminé par $V(i,j)$, il se dirige vers la cellule du dessus
- Si l'élément $S(i,j)$ est déterminé par $W(i,j)$, il se dirige vers la cellule de gauche. Ces deux éléments introduisent alors des gap dans la solution
- Si l'élément $S(i,j)$ est déterminé par $S(i-1, j-1)$ + son score de substitution, il se dirige en diagonale en haut à gauche de lui.

$$\begin{bmatrix} S(i-1, j-1) + t(i,j) & V(i,j) \\ W(i,j) & S(i,j) \end{bmatrix}$$

```
In [16]: 1 def is_previous(self, mat, i, j):
2         """ Regarde si l'élément ij résulte de l'élément en diagonale,
3             en haut ou a gauche
4         """
5
6         res = False
7         if mat == "v":
8             if self.V.get_score(i,j) == self.S.get_score(i,j):
9                 res = True
10
11        elif mat == "w":
12            if self.W.get_score(i,j) == self.S.get_score(i,j):
13                res = True
14
15        elif mat == "s":
16            letters_ij = self.S.get_letters(i,j) # 'AB' par exemple
17            t_ij = self.t.get_acid_score(letters_ij[0], letters_ij[1]) # t(i,j)
18
19            if self.S.get_score(i-1, j-1) + t_ij == self.S.get_score(i,j):
20                res = True
21
22        return res
23
24 Alignment.is_previous = is_previous
25
```

Alignement local

L'alignement local va utiliser l'algorithme de Smith-Waterman. Nous voulons trouver les morceaux de séquences similaires. Cet algorithme utilise également la programmation dynamique et ressemble à l'algorithme de Needleman-Wunsch.


```

In [17]: 1 class LocalAlignment(Alignment):
2         """
3         Classe qui va trouver un aligement de séquences d'acides aminées
4         une pénalité affine et une méthode locale
5         """
6
7         def __init__(self, l, I, E, mat_file, seq1, seq2, p):
8             Alignment.__init__(self, I, E, mat_file, seq1, seq2, p)
9             self.l = l
10
11             self.S.init_S_local() # On initialise la matrice de scoring
12             self.S.set_letters_seq("-"+self.seq1.get_acids(), "-"+self.seq2.get_acids())
13
14             self.zeros = []
15             self.found = False
16
17         def update_s_local(self, i, j, v_ij, w_ij):
18             """ détermine la valeur de S en fonction de V et W et de t """
19
20             letters_ij = self.S.get_letters(i,j) # 'AB' par exemple
21             t_ij = self.t.get_acid_score(letters_ij[0], letters_ij[1])
22             s_ij = max( self.S.get_score(i-1,j-1) + t_ij, v_ij, w_ij, 0 )
23             self.S.set_score(i,j, s_ij)
24
25

```

Tout d'abord, nous initialisons les matrices V, W et S. V et W sont de la même forme que ceux utilisés pour l'alignement global:

$$V = \begin{bmatrix} -Inf & -Inf & \dots \\ 0 & & \\ 0 & & \\ \dots & & \end{bmatrix} \quad W = \begin{bmatrix} 0 & 0 & \dots \\ -Inf & & \\ -Inf & & \\ \dots & & \end{bmatrix}$$

La matrice de scoring ne contiendra que des valeurs positives ou nulles. Si nous obtenons un score négatif, celui-ci devient 0. De ce fait, la matrice est initialisée de cette manière:

$$S = \begin{bmatrix} 0 & 0 & 0 & \dots \\ 0 & & & \\ 0 & & & \\ \dots & & & \end{bmatrix}$$

La méthode ci-dessous réalise cette initialisation:

```

In [18]: 1 def init_S_local(self):
2         """ Initialise 1ere ligne et 1ere colonne de S pour la méthode locale
3
4         for i in range(self.m):
5             newline = []
6             if i == 0:
7                 newline += [0]*self.n
8             else:
9                 newline += [0] + (self.n-1)*[""]
10
11             self.set_lign(newline)
12
13 MatScoring.init S local = init S local

```

Une fois les matrices V, W et S initialisée, nous les calculons à l'aide de ces relations:

$$V(i,j) = \max \begin{cases} S(i-1,j) - I \\ V(i-1,j) - E \end{cases}$$

$$W(i,j) = \max \begin{cases} S(i,j-1) - I \\ W(i,j-1) - E \end{cases}$$

$$S(i,j) = \max \begin{cases} S(i-1,j-1) + \text{MatSubstitution}(i,j) \\ V(i,j) \\ W(i,j) \\ 0 \end{cases}$$

Ceci est réalisé par la méthode ci-dessous:

```
In [19]: 1 def compute_scoring(self, start_i, start_j):
2         """ ReCalcule la matrice de
3         scoring pour les alignements locaux avec pénalité affine
4         après avoir trouvé un alignement local
5         """
6
7         # ===== CREATION SCORING =====
8
9         for i in range(start_i, self.m):
10            for j in range(start_j, self.n):
11                if (i,j) not in self.zeros:
12                    v_ij = self.get_v(i, j) # V(i,j)
13                    self.V.set_score(i,j, v_ij)
14                    w_ij = self.get_w(i, j) # W(i,j)
15                    self.W.set_score(i,j, w_ij)
16
17                    self.update_s_local(i, j, v_ij, w_ij)
18
19            if self.p == 2:
20                print("la matrice de Scoring: ")
21                self.S.panda()
22
23 LocalAlignment.compute_scoring = compute_scoring
```

Une fois la matrice de scoring, ainsi que les matrices de sauvegarde V et W sont calculées, nous cherchons la cellule dont le score est maximal. Nous utiliserons la méthode prédéfinie ci-dessous:

```
In [20]: 1 def get_max(self):
2         """ Renvoie la position de l'élément maximal de la matrice """
3
4         maxi = -float("inf") # - inf
5         i, j = 0, 0
6         for line in range(self.m):
7             current_max = max(self.mat[line])
8             if current_max > maxi:
9                 maxi = current_max
10                i = line
11                j = self.mat[line].index(current_max)
12
13         return (i,j)
14
15 Matrix.get_max = get_max
```

Depuis ce maximum, nous réalisons un backtracking qui ne recherche qu'une solution. Celui-ci part donc de la cellule dont le score est maximal et remonte selon les mêmes conditions que l'alignement global.

- Si contre, la cellule $S(i,j)$ a été déterminé par le 0 lors du calcul du maximum (voir relations ci-dessus), nous nous arrêtons.
- Si il ne rencontre pas de 0, il s'arrête lorsqu'il se trouve sur la première ligne ou sur la première colonne de la matrice score.

Cette approche "bottom-up" est mise en oeuvre par les méthodes ci-dessous:

```
In [21]: 1 def sol_found(self, i,j):
2         """ Détermine si on a fini un alignement local """
3
4         return (i == 0 or j == 0 or self.S.get_score(i,j) == 0)
5
6     LocalAlignment.sol_found = sol_found
7
8     def bottom_up(self, i,j):
9         """ Remonte la matrice de scoring a partir du max de la matrice
10             de scoring jusqu'à un élément de la 1ere ligne ou 1ere colonne
11             ou un élément 0
12         """
13         self.current_sol.append((i,j))
14         while not self.sol_found(i,j):
15             if self.is_previous("v", i, j): # haut
16                 i-=1
17             elif self.is_previous("w", i, j): # gauche
18                 j-=1
19             elif self.is_previous("s", i, j): # diagonale
20                 i-=1
21                 j-=1
22             if self.S.get_score(i,j) != 0:
23                 self.current_sol.append((i,j))
24
25         self.all_solutions.append(self.current_sol)
26         self.zeros += self.current_sol
27         self.current_sol = []
28
29     LocalAlignment.bottom up = bottom up
```

Après avoir enregistré la solution (le chemin), nous mettons les éléments sur lesquels nous sommes passés à 0.

```
In [22]: 1 def set_zero(self, positions):
2         """ Met à 0 les éléments donnés en paramètre dans la matrice """
3
4         for pos in positions:
5             self.mat[pos[0]][pos[1]] = 0
6
7     Matrix.set zero = set zero
```

Enfin, étant donné que nous souhaitons connaître un certain nombre de solution, il nous faut répéter ce processus. Pour résumer, l'algorithme Smith-Waterman ci-dessous va réaliser les actions suivantes: Tant que l'on n'a pas trouvé k solutions:

- Trouver la cellule ayant le score maximal de la matrice de scoring préalablement calculée
- Procéder à un "bottom up" pour trouver le chemin correspondant et en s'arrêtant soit en ayant rencontré un 0, soit en étant sur la première ligne ou la première colonne
- Mettre à 0 ce nouveau chemin dans la matrice de scoring et dans les matrices V et W
- Recalculer les 3 matrices S, V et W en tenant compte des 0 mis en place précédemment.

```
In [23]: 1 def Smith_Waterman(self):
2         self.compute_scoring(1,1)
3
4         for i in range(self.l):
5             current_max = self.S.get_max()
6             self.bottom_up(current_max[0], current_max[1])
7             if i == self.l-1: # Si on a trouvé l solutions
8                 break
9             # une fois le backtrack fini, on met a 0 l'alignement
10            self.S.set_zero(self.all_solutions[-1])
11            self.V.set_zero(self.all_solutions[-1])
12            self.W.set_zero(self.all_solutions[-1])
13            if self.p == 2:
14                print("mise a 0: ")
15                self.S.panda()
16
17                print("recalcul de la matrice: ")
18                self.compute_scoring(self.all_solutions[-1][-1][0], \
19                                    self.all_solutions[-1][-1][1])
20
21            R = Result(self.S, self.t, self.all_solutions, self.p)
22            return R.compute_result()
23
24 LocalAlignment.Smith Waterman = Smith Waterman
```

Résultats & Discussion

ADT Resultat

Cette classe nous permet de mettre en place la solution calculée grâce aux algorithmes. Elle va d'abord rassembler les paires d'acides aminés en instaurant des gap selon les conditions décrites plus haut (sur base des matrices V et W).

```

In [24]: 1 class Result:
2         """ Classe représentant un objet Résultat dans laquelle on va aligner
3           2 séquences selon la matrice de scoring obtenue
4         """
5
6         def __init__(self, S, t, all_sol, p):
7             self.S = S # matrice de scoring
8             self.t = t # matrice de substitution
9             self.all_solutions = all_sol
10            self.gap = "-"
11            self.p = p # Pour savoir si on veut print le résultat ou non
12
13        def compute_result(self):
14            """ trouve les correspondances entre lettres de l'alignement """
15
16            scores_sim = [] # stocke le score de similarité des résultats
17
18            for sol in self.all_solutions:
19                used = {} # dictionnaire avec indices déjà utilisés
20                        # clé = colonnes, valeur = lignes
21                res = []
22                for pos in range(len(sol)-1, -1, -1):
23
24                    letters = self.S.get_letters(sol[pos][0], sol[pos][1])
25                    score = self.t.get_acid_score(letters[0], letters[1])
26                    if sol[pos][0] in used.values():
27                        letters = self.gap + letters[1] # '-B' par ex
28                    if sol[pos][1] in used.keys():
29                        letters = letters[0] + self.gap # 'B-' par ex
30
31                    res.append((letters, score)) # On ajoute son score aussi
32
33                    used[sol[pos][1]] = sol[pos][0]
34
35                #self.bind(res)
36                scores_sim.append(self.bind(res))
37
38            return scores_sim

```

Une méthode de la classe résulte nous permet de construire les similitudes entre les 2 séquences.

- Si 2 acides aminés sont identiques, on ajoute "."
- Si 2 acides aminés sont différents mais similaires d'un point de vue de l'évolution biologique, c'est-à-dire que leur score dans la matrice de substitution est positive, on ajoute un "."
- Si il y a un gap (représenté par un tiret), ou si le score est négatif, on ajoute un espace " "

```

In [25]: 1 def bind(self, seq):
          2     """ crée la liaison entre les 2 séquences (similarité, identité, ..)
          3         ainsi que les scores et pourcentages
          4     """
          5
          6     # print(seq)
          7     seq1, seq2, links = "", "", ""
          8     for i in range(len(seq)):
          9         # print("seq: ", seq[i][0], " score: ", seq[i][1])
         10         seq1 += seq[i][0][1] # 2eme lettre
         11         seq2 += seq[i][0][0] # 1ere lettre
         12
         13         if seq1[i] == self.gap or seq2[i] == self.gap:
         14             links += " " # Pas de correspondance
         15         else:
         16             if seq1[i] == seq2[i]:
         17                 links += ":" # Identiques
         18             elif seq[i][1] >= 0: # Si le score est positif => similaire
         19                 links += "."
         20             else:
         21                 links += " "
         22
         23         similarity = ((links.count(".") + links.count(":")) / len(links)) * 100 #
         24         identity = (links.count(":") / len(links)) * 100 # pourcentage d'identi
         25
         26         if self.p == 1 or self.p == 2:
         27             self.print_result(seq1, seq2, links, similarity, identity)
         28
         29     return similarity
         30
         31 Result.bind = bind

```

Une fois que nous avons obtenu les paires d'acides aminés et leurs liens respectifs, nous pouvons l'afficher de manière élégante à l'aide de la méthode ci-dessous.

```

In [26]: 1 def print_result(self, seq1, seq2, links, similarity, identity):
          2     """ Imprime de manière jolie les 2 séquences alignées
          3         seq est de la forme: [('AB', scoreAB), .. ]
          4     """
          5
          6     print("\n>>> Alignement: ")
          7
          8     i, j = 0, 0
          9     while i < (len(seq1) // 60):
         10         print(seq1[j:j+60] + "\n" + links[j:j+60] + "\n" + seq2[j:j+60] + "\n\n")
         11         i += 1
         12         j += 60
         13     end = len(seq1) - len(seq2) % 60
         14     print(seq1[end:] + "\n" + links[end:] + "\n" + seq2[end:] + "\n\n")
         15
         16     print("==> similarité: {0} %".format(similarity))
         17     print("==> Identité: {0} %".format(identity))
         18
         19 Result.print_result = print_result

```

Global

Afin de comparer les résultats obtenus, nous utiliserons le site "LALIGN" qui procède également à des alignements de séquences. Nous utiliserons également le site "UNIPROT" qui répertorie un grand nombre de protéines.

Exemple simple

Voici un exemple simple d'alignement global de 2 petites séquences avec sa matrice de scoring:

```
In [27]: 1 seq1 = Sequence("Séquence 1 de l'exemple", "MGGETFA")
          2 seq2 = Sequence("Séquence 2 de l'exemple", "GGVTTF")
          3 G0 = GlobalAlignment(2, 12, 2, "blosum62.txt", seq1, seq2, 2)
          4 G0.Needleman Wunsch()
```

Séquence 1 de longueur 8:

Séquence 1 de l'exemple
MGGETFA

Séquence 2 de longueur 7:

Séquence 2 de l'exemple
GGVTTF

matrice de substitution utilisée: blosum62.txt

Pénalité de gap affine: I = 12 | E = 2

*	-	M	G	G	E	T	F	A
-	0	-12	-14	-16	-18	-20	-22	-24
G	-12	-3	-6	-8	-18	-20	-23	-22
G	-14	-15	3	0	-10	-13	-15	-17
V	-16	-13	-9	0	-2	-10	-14	-15
T	-18	-17	-11	-11	-1	3	-9	-11
T	-20	-19	-13	-13	-12	4	1	-9
F	-22	-20	-15	-16	-15	-8	10	-1

>>> Alignement:

MGGETFA

:::

-GGVTTF

=> similarité: 42.857142857142854 %

=> Identité: 42.857142857142854 %

Out[27]: [42.857142857142854]

Comme nous pouvons le voir, cet exemple coïncide avec celui présenté dans les slides du projet.

Séquences du fichier BRD-sequence

Nous testons maintenant d'aligner certaines séquences d'acides aminés contenues dans le fichier BRD-sequence.fasta

Données de UNIPROT :

organisme de toutes les séquences: Homo sapiens (Human)

- Séquence 1 (O60885) :
 - protéine = Bromodomain-containing protein 4,
 - gene = BRD4,
- Séquence 2 (O60885) : idem que séquence 1
- Séquence 3 (P21675) :
 - protéine = Transcription initiation factor TFIID subunit 1,
 - gene = TAF1
- Séquence 4 (P21675) : idem que séquence 3
- Séquence 5 (Q92831) :
 - protéine = Histone acetyltransferase KAT2B,
 - gene = KAT2B
- Séquence 6 (Q86U86) :
 - protéine = Protein polybromo-1
 - gene = PBRM1
- Séquence 7 (Q9NPI1) :
 - protéine = Bromodomain-containing protein 7,
 - gene = BRD7

```
In [28]: 1 P = ParserSequence("BRD-sequence.fasta")
          2 P.parse()
```

Séquences les plus similaires parmi les séquences du fichier

Nous calculons toutes les combinaisons de séquences du fichier "BRD-protein.fasta" et nous en sortons les 2 ayant le plus de similitudes.

==> Identité: 50.70422535211267 %

```

In [30]: 1 G2 = GlobalAlignment(1, 12, 2, "blosum62.txt", P.get_seq(0), P.get_seq(0)
          2 G2.Needleman Wunsch()

          Séquence 1 de longueur 74:
          >sp|060885|75-147

          WKHQFAWPFQQPVDAVKLNLPDYYKIIKTPMDMGTIKKRENNYYWNAQECIQDFNTMFTNCYIYNKPGDDIV

          Séquence 2 de longueur 74:
          >sp|060885|75-147

          WKHQFAWPFQQPVDAVKLNLPDYYKIIKTPMDMGTIKKRENNYYWNAQECIQDFNTMFTNCYIYNKPGDDIV

          matrice de substitution utilisée: blosum62.txt
          Pénalité de gap affine: I = 12 | E = 2

          >>> Alignement:
          WKHQFAWPFQQPVDAVKLNLPDYYKIIKTPMDMGTIKKRENNYYWNAQECIQDFNTMFT
          .....
          WKHQFAWPFQQPVDAVKLNLPDYYKIIKTPMDMGTIKKRENNYYWNAQECIQDFNTMFT

          NCYIYNKPGDDIV
          .....
          NCYIYNKPGDDIV

          ==> similarité: 100.0 %
          ==> Identité: 100.0 %

```

Out[30]: [100.0]

Comme prévu, lorsque nous procédons à l'alignement de 2 séquences identiques, ceux-ci ont un pourcentage de similarité de 100 %.

Séquences 1 et 2 du fichier BRD-sequences

```
In [31]: 1 G3 = GlobalAlignment(1, 12, 2, "blosum62.txt", P.get_seq(0), P.get_seq(1)
          2 G3.Needleman_Wunsch()
          3

          Séquence 1 de longueur 74:
          >sp|060885|75-147

          WKHQFAWPFQPPVDAVKLNLDPYYKIIKTPMDMGTIKKRLNNYYWNAQECIQDFNTMFTNCYIYNKPGDDIV

          Séquence 2 de longueur 74:
          >sp|060885|368-440

          KHAAYAWPFYKPPVDVEALGLHDYCDIIKHPMDMSTIKSKLEAREYRDAQEFGADVRLMFSNCYKYNPPDHEVV

          matrice de substitution utilisée: blosum62.txt
          Pénalité de gap affine: I = 12 | E = 2

          >>> Alignement:
          WKHQFAWPFQPPVDAVKLNLDPYYKIIKTPMDMGTIKKRLNNYYWNAQECIQDFNTMFT
          .:..: .:..: .:..: .:..: .:..: .:..: .:..: .:..: .:..: .:..: .:..:
          KHAAYAWPFYKPPVDVEALGLHDYCDIIKHPMDMSTIKSKLEAREYRDAQEFGADVRLMFS

          NCYIYNKPGDDIV
          :..: .:..: .:..: .:..:
          NCYKYNPPDHEVV

          ==> similarité: 68.4931506849315 %
          ==> Identité: 50.68493150684932 %
```

```
Out[31]: [68.4931506849315]
```

```
Algorithm: Global/Global affine Needleman-Wunsch (SSE2, Michael Farrar 2010) (6.0 April 2007)
Parameters: data/blosum62.mat matrix (11:-4), open/ext: -12/-2
Scan time: 0.000

The best scores are:
unknown 73 bp                                n-w bits E(1)
                                         ( 73) 182 48.3 2.3e-80

>>unknown 73 bp                                (73 aa)
n-w opt: 182 Z-score: 239.5 bits: 48.3 E(1): 2.3e-80
global/global (N-W) score: 182; 50.7% identity (68.5% similar) in 73 aa overlap (1-73:1-73)

      10      20      30      40      50      60
unknown WKHQFAWPFQPPVDAVKLNLDPYYKIIKTPMDMGTIKKRLNNYYWNAQECIQDFNTMFT
      .:..: .:..: .:..: .:..: .:..: .:..: .:..: .:..: .:..: .:..:
unknown KHAAYAWPFYKPPVDVEALGLHDYCDIIKHPMDMSTIKSKLEAREYRDAQEFGADVRLMFS
      10      20      30      40      50      60

      70
unknown NCYIYNKPGDDIV
      :..: .:..: .:..: .:..:
unknown NCYKYNPPDHEVV
      70
```

Nous pouvons constater que les résultats obtenus avec notre algorithme sont identiques à ceux de LALIGN. En effet, nous avons utilisé la même matrice de substitution "BLOSUM62" ainsi que des pénalités de gap identiques (I = 12 et E = 2). L'algorithme utilisé par LALIGN est Needleman-Wunsch.

De plus, le taux de similitudes (68.5%) et le taux d'identité (50.7%) des 2 sont identiques.

En regardant les données de UNIPROT, nous nous apercevons que les 2 séquences appartiennent à la même famille: un bromodomaine contenant une protéine 4. Le gène est également identique: BRD4 et les 2 sont issus d'un organisme Humain.

Séquence 1 et 3 du fichier BRD-sequences

```
In [32]: 1 G4 = GlobalAlignment(1, 12, 2, "blosum62.txt", P.get_seq(0), P.get_seq(2)
2 G4.Needleman Wunsch()

Séquence 1 de longueur 74:
>sp|060885|75-147

WKHQFAWPFQPPVDAVKLNLPDYYKIIKTPMDMGTIKKRENNYYWNAQECIQDFNTMFTNCYIYNKPGDDIV

Séquence 2 de longueur 72:
>sp|P21675|1397-1467

RDLPNITYPFHTPVNAKVVKDYYKIIITRPMDLQTLRENVKRKLYPSREEFRELIVKNSATYNGPKHSLT

matrice de substitution utilisée: blosum62.txt
Pénalité de gap affine: I = 12 | E = 2

>>> Alignement:
WKHQFAWPFQPPVDAVKLNLPDYYKIIKTPMDMGTIKKRENNYYWNAQECIQDFNTMFT
..... :::: .. ::::: :::: :..... : . : . . . .
RDLPNITYPFHTPVNA-KVVK-DYYKIIITRPMDLQTLRENVKRKLYPSREEFRELIVK

NCYIYNKPGDDIV
:  :: : ...
NSATYNGPKHSLT

==> similarité: 61.64383561643836 %
==> Identité: 30.136986301369863 %
```

```
Out[32]: [61.64383561643836]
```

```
Algorithm: Global/Global affine Needleman-Wunsch (SSE2, Michael Farrar 2010) (6.0 April 2007)
Parameters: data/blosum62.mat matrix (11:-4), open/ext: -12/-2
Scan time: 0.000

The best scores are:
unknown 71 bp
n-w bits E(1)
( 71) 82 33.3 8.6e-28

>>unknown 71 bp
n-w opt: 82 Z-score: 158.6 bits: 33.3 E(1): 8.6e-28
global/global (N-W) score: 82; 28.8% identity (61.6% similar) in 73 aa overlap (1-73:1-71)

      10      20      30      40      50      60
unknown WKHQFAWPFQPPVDAVKLNLPDYYKIIKTPMDMGTIKKRENNYYWNAQECIQDFNTMFT
      :::: :::: .. ::::: :::: :..... : . : . . . .
unknown RDLPNITYPFHTPVNAKVVK--DYYKIIITRPMDLQTLRENVKRKLYPSREEFRELIVK
      10      20      30      40      50

      70
unknown NCYIYNKPGDDIV
      :  :: : ...
unknown NSATYNGPKHSLT
      60      70
```

Avec ces 2 séquences (1er et 3ème du fichier BRD-sequences), nous obtenons le même résultat que LALIGN à peu de chose près. En effet, la sous séquence "KVVK" est entourée par 1 gap de chaque côté dans notre cas, tandis que d'après LALIGN, à la position 20, la sous-séquence "KVVK" est suivie de 2 gap. Cela peut être dû à la pénalité de gap affine choisie.

En regardant les données de UNIPROT, nous avons 2 séquences différentes:

- La première séquence est un "Bromodomaine contenant une protéine 4", et de gène BRD4
- La deuxième séquence est un "Transcription initiation factor TFIID subunit 1", et de gène TAF1

Pourtant, ils ont un taux de similarités de 61,6 %. nous poserons l'hypothèse que ces 2 séquences d'un organisme de type Humain sont assez similaires malgré qu'ils ne soient pas de la même famille.

Local

Afin de tester les solutions des alignements locaux, nous utiliserons également les sites "LALIGN" pour les alignements et "UNIPROT" pour les données liées au type de protéines.

Exemple simple

Voici un alignement local sur un exemple simple avec la matrice de substitution BLOSUM62 et une pénalité affine ($I = 4$, $E = 4$) :

```
In [33]: 1 P2 = ParserSequence("protein-sequences.fasta")
          2 P2.parse()
          3
          4 seq1 = Sequence("Séquence 1 de l'exemple", "ISALIGNED")
          5 seq2 = Sequence("Séquence 2 de l'exemple", "THISLINE")
          6 L0 = LocalAlignment(2, 4, 4, "blosum62.txt", seq1, seq2, 2)
          7 L0.Smith Waterman()
```

Séquence 1 de longueur 10:

Séquence 1 de l'exemple

ISALIGNED

Séquence 2 de longueur 9:

Séquence 2 de l'exemple

THISLINE

matrice de substitution utilisée: blosum62.txt

Pénalité de gap affine: I = 4 | E = 4

la matrice de Scoring:

*	-	I	S	A	L	I	G	N	E	D
-	0	0	0	0	0	0	0	0	0	0
T	0	0	1	0	0	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0
I	0	4	0	0	2	4	0	0	0	0
S	0	0	8	4	0	0	4	1	0	0
L	0	2	4	7	8	4	0	1	0	0
I	0	4	0	3	9	12	8	4	0	0
N	0	0	5	1	5	8	12	14	10	6
E	0	0	1	4	1	4	8	12	19	15

mise a 0:

*	-	I	S	A	L	I	G	N	E	D
-	0	0	0	0	0	0	0	0	0	0
T	0	0	1	0	0	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0
I	0	0	0	0	2	4	0	0	0	0
S	0	0	0	0	0	0	4	1	0	0
L	0	2	4	7	0	4	0	1	0	0
I	0	4	0	3	9	0	0	4	0	0
N	0	0	5	1	5	8	12	0	10	6
E	0	0	1	4	1	4	8	12	0	15

recalcul de la matrice:

la matrice de Scoring:

*	-	I	S	A	L	I	G	N	E	D
-	0	0	0	0	0	0	0	0	0	0
T	0	0	1	0	0	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0
I	0	0	0	0	2	4	0	0	0	0
S	0	0	0	0	0	0	4	1	0	0
L	0	2	0	0	0	2	0	1	0	0
I	0	4	0	0	2	0	0	0	0	0
N	0	0	5	1	0	0	0	0	0	1
E	0	0	1	4	0	0	0	0	0	2

>>> Alignement:

ISALIGNE

:: :: ::

IS-LI-NE

==> similarité: 75.0 %

==> Identité: 75.0 %

>>> Alignement:

IS

..

TM

Comme nous pouvons le voir, cet exemple coïncide avec celui présenté dans les slides du projet.

Séquences du fichier protein-sequences

Nous allons maintenant tester l'alignement local sur différentes séquences du fichier protein-sequence.fasta

Données de UNIPROT:

Organisme de toutes les séquences: Homo sapiens (Human)

- Séquence 1 (Q9H8M2) :
 - protéine = Bromodomain-containing protein 9
 - gene = BRD9
- Séquence 2 (O60885) :
 - protéine = Bromodomain-containing protein 4
 - gene = BRD4
- Séquence 3 (P21675) :
 - protéine = Transcription initiation factor TFIID subunit 1
 - gene = TAF1
- Séquence 4 (Q86U86) :
 - protéine = Protein polybromo-1
 - gene = PBRM1
- Séquence 5 (Q9NPI1):
 - protéine = Bromodomain-containing protein 7
 - gene = BRD7

Séquences identiques

Nous testons l'alignement local sur 2 séquences identiques :

```
In [34]: 1 L1 = LocalAlignment(1, 12, 2, "blosum62.txt", P2.get_seq(0), P2.get_seq(6)
2 L1.Smith Waterman()

Séquence 1 de longueur 598:
>sp|Q9H8M2|BRD9_HUMAN Bromodomain-containing protein 9 OS=Homo sapiens OX
=9606 GN=BRD9 PE=1 SV=2

MGKKHKHHKA EWRSSYEDYADKPLEKPLKLVLKVGGSEVTELSGSGHDSSYYDDRS DHERERHKEKKKKKKKK
SEKEKHL DDEERRKRKEEKKRKREREHC DTEGEADD FPGKKVEVEPPDRPV RACRTQPAENESTPIQLLE
HFLRQLQRKD PHGFFAFPVTD AIAPGYSMIIKH PMDFGT MKDKIVANEYKS VTEFKADF KLMCDNAMTYNRPD
TVYYKLAKKIL HAGFKMMSQA ALLGNEDT AVEEPVPEVP VPQVET AKSKKPSREVIS CMFEPEG NACSLTD
STAEHVLA LVEHAAD EARDRI NRFLP GKGMYL KRNGDGS LLYSVVNTA EPDADEEETH PVDSLSSL SKLLP
GFTTLG FKDERRN KVTF LSSATTAL SMQNNSVF GD LKSD EMELLY SAYGDETGVQC ALSLQEFVK DAGSYSKK
VVDDL LDQITGG DHSRTL FQLK QRRNVPM KPDPDE AKVGDT LGDSSSSSV LEFMSMK SYPDVSVD ISMLSSL GKV
KKELD PDSSHNLN DETTKLL QDLHEAQ AERGGR PSSNL SSLNAS ERDQHHLGSP SRLSVGEQP DVTHDPYE
FLQSPE PAASAKT

Séquence 2 de longueur 598:
>sp|Q9H8M2|BRD9_HUMAN Bromodomain-containing protein 9 OS=Homo sapiens OX
=9606 GN=BRD9 PE=1 SV=2

MGKKHKHHKA EWRSSYEDYADKPLEKPLKLVLKVGGSEVTELSGSGHDSSYYDDRS DHERERHKEKKKKKKKK
SEKEKHL DDEERRKRKEEKKRKREREHC DTEGEADD FPGKKVEVEPPDRPV RACRTQPAENESTPIQLLE
HFLRQLQRKD PHGFFAFPVTD AIAPGYSMIIKH PMDFGT MKDKIVANEYKS VTEFKADF KLMCDNAMTYNRPD
TVYYKLAKKIL HAGFKMMSQA ALLGNEDT AVEEPVPEVP VPQVET AKSKKPSREVIS CMFEPEG NACSLTD
STAEHVLA LVEHAAD EARDRI NRFLP GKGMYL KRNGDGS LLYSVVNTA EPDADEEETH PVDSLSSL SKLLP
GFTTLG FKDERRN KVTF LSSATTAL SMQNNSVF GD LKSD EMELLY SAYGDETGVQC ALSLQEFVK DAGSYSKK
VVDDL LDQITGG DHSRTL FQLK QRRNVPM KPDPDE AKVGDT LGDSSSSSV LEFMSMK SYPDVSVD ISMLSSL GKV
KKELD PDSSHNLN DETTKLL QDLHEAQ AERGGR PSSNL SSLNAS ERDQHHLGSP SRLSVGEQP DVTHDPYE
FLQSPE PAASAKT

matrice de substitution utilisée: blosum62.txt
Pénalité de gap affine: I = 12 | E = 2

>>> Alignement:
MGKKHKHHKA EWRSSYEDYADKPLEKPLKLVLKVGGSEVTELSGSGHDSSYYDDRS DHER
::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::
MGKKHKHHKA EWRSSYEDYADKPLEKPLKLVLKVGGSEVTELSGSGHDSSYYDDRS DHER

ERHKEKKKKKKK SEKEKHL DDEERRKRKEEKKRKREREHC DTEGEADD FPGKKVEVEP
::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::
ERHKEKKKKKKK SEKEKHL DDEERRKRKEEKKRKREREHC DTEGEADD FPGKKVEVEP

PPDRPV RACRTQPAENESTPIQLLE HFLRQLQRKD PHGFFAFPVTD AIAPGYSMIIKH P
::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::
PPDRPV RACRTQPAENESTPIQLLE HFLRQLQRKD PHGFFAFPVTD AIAPGYSMIIKH P

MDFGTM KDKIVANEYKS VTEFKADF KLMCDNAMTYNRPD TVYYKLAKKIL HAGFKMMSKQ
::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::
MDFGTM KDKIVANEYKS VTEFKADF KLMCDNAMTYNRPD TVYYKLAKKIL HAGFKMMSKQ

AALLGNEDTAVEEPVPEVP VPQVET AKSKKPSREVIS CMFEPEG NACSLTD STAEHVLA
::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::
AALLGNEDTAVEEPVPEVP VPQVET AKSKKPSREVIS CMFEPEG NACSLTD STAEHVLA

ALVEHAAD EARDRI NRFLP GKGMYL KRNGDGS LLYSVVNTA EPDADEEETH PVDSLSSL
::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::
ALVEHAAD EARDRI NRFLP GKGMYL KRNGDGS LLYSVVNTA EPDADEEETH PVDSLSSL

SKLLPG FTTLG FKDERRN KVTF LSSATTAL SMQNNSVF GD LKSD EMELLY SAYGDETGVQ
::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::
SKLLPG FTTLG FKDERRN KVTF LSSATTAL SMQNNSVF GD LKSD EMELLY SAYGDETGVQ
```


Comme prévu, le taux de similarité est maximal étant donné qu'ils sont identiques.

Séquences les plus similaires parmi le fichier protein-sequences

Nous essayons à présent de déterminer les séquences du fichier protein-sequences.fasta les plus similaires avec l'alignement local :

```
In [35]: 1 def get_best_localalignments():
2         best_score = [0, 0, 0] # les 2 derniers 0 accueilleront les numéros d'indices
3                                     # le score est maximal
4
5         for i in range(5):
6             for j in range(5):
7                 if i != j: # on ne compare pas une séquence avec elle-même
8                     #print(i, " ", j)
9                     L_tmp = LocalAlignment(1, 12, 2, "blosum62.txt", P2.get_seq(1))
10                    tmp_score = L_tmp.Smith_Waterman()
11                    if tmp_score[0] > best_score[0]:
12                        best_score[0] = tmp_score[0]
13                        best_score[1] = i
14                        best_score[2] = j
15            return best_score
16 best = get_best_localalignments()
17
18 print("le meilleur score est celui-ci: ", best[0])
19 L2 = LocalAlignment(1, 12, 2, "blosum62.txt", P2.get_seq(best[1]), P2.get_seq(best[2]))
20 L2 = L2.Smith_Waterman()
```

le meilleur score est celui-ci: 71.42857142857143

Séquence 1 de longueur 1873:

>sp|P21675|TAF1_HUMAN Transcription initiation factor TFIID subunit 1 OS=Homo sapiens OX=9606 GN=TAF1 PE=1 SV=2

MGP GCDLLLR TAATITAAAIMS DTDSD EDSAGGGP FSLAGFLFG NINGAGQLEGESVLDDECKKHLA GLGALG
 LGSLITELTANEELTGTDGALVNDEGWVRSTEDAVDYS DINEVAEDESRRYQQTMSGLQPLCHSDYDEDDYDA
 DCEDIDCKLMPPPPPPGPMKKDKDQDSITGEKVDFSSSDSESEMGPQEATQAESEDKGLTLPAGIMQHDA
 TKLLPSVTELFPEFRPGKVLRLFLFGPGKNVPSVWRSARRKRKKKHREL IQEEIQEVECSVESEVSQKSLW
 NYDYAPPPPEQCLSDDEITMMA PVESKFSQSTGDI DKVTDTKPRVAEWRYGPARLWYDMLGVPEDGSGFDYG
 FKLRKTEHEPVIKSRMIEEFRKLEENNGTDLLADENFLMVTQLHWEDDI IWDGEDVKHKGTKPQRASLAGWLP
 SSMTRNAMAYNVQQGFAATLDDDKPWYSIFPIDNEDLVYGRWEDNIIWDAQAMPRLLEPPVLTLPNDENLIL
 EIPDEKEEATSNSPSKESKKESSLKKSRI LLGKTGVIKEEPQONMSQPEVKDPWNLSNDEYYYPKQQGLRGTF
 GGNIIQHSIPAVELRQPFPTHMGPIKLRQFHRPPLKKYSFGALSQPGPHSVQPLLKHKKKAKMREQERQAS
 GGGEMFFMRTPQDLTGKDGDLILA EYSEENGPLMMQVGMATKIKNYYKRKPGKDPGAPDCKYGETVYCHTSPF
 LGSLHPGQLLQAFENNLFRAPIYLHKMPETDFLIIRTRQGYI IRELVDIFVVGQQCPLEFVPGPN SKRANTHI
 RDFLQVFIYRLFWKSKDRPRRIRMEDIKKAFPSHSESSIRKRLKLCADFKRTGMDSNWVWLKSDFRLPTEEEI
 RAMVSPEQCCAYYSMIAAEQRLKDAGYGEKSFFAPEEENEEDFQMKIDDEVRTAPWNTTTRAFIAAMKGKCLLE
 VTGVADPTGCGEGFSYVKIPNKPTQKDDKEPQPVKKTVTGTADLRRLSLKNAKQLLRKFGVP EEEIKKLSR
 WEVIDVVRTMSTEQARSGEGPMSK FARGSRFSVAEHQERYKEECQRIFDLQNKVLSSTEVLTSDT DSSSAEDS
 DFEEMGKN IENMLQNKKTSSQLSRERE EQRKELQRMLLAAGSAASGNNHRDDDTASVTS LNSSATGRCLKIY
 RTFRDEEGKEYVRCETVRKPAVIDAYVRITTKDEEFIRKFALFDEQHREEMRKERRRIQEQLRRLKRNQEKE
 KLKGPPEKKPKMKKERPD LKLKCGACGAIGHMRTNKFCLPYQT NAPPSPNPVAMTEE EEEEELEKTVIHNDNEE
 LIKVEGTKIVLGKQLIESADEVRRKSLVLKFKPQQLP PPKKKRRVGTTVHCDYLNRP HKSIIHRRRTDPMVLTSS
 ILESIIINDMRDL PNTYPFHTPVNAKVVDYKIIITRPMDLQTLRENVKRRLYPSREEF REHLELIVKNSATYN
 GPKHSLTQISQSM L DCLDEKLKEKEDKLARLEKAINPL LDDDDQVAFS FILDNI VTQKMMAVPDSWPFHHPVN
 KKFVPDYKYVIVNPM DLETIRKNISKHKYQSRESFLDDVNLILANSVKYNGPESQYTKTAQEI VNVCYQTLTE
 YDEHLTQLEKDICTAKEAAL EEALES LDPMTPGPYTPQPPDLYDTNTSLSMSRDASVFQDES NMSVLDIPSA
 TPEKQVTQEGEDGDGLADEEEGT VQPPQASVLYEDLLMSEGEDDEEDAGSDEEGDNPFS AIQLSESGSDSDV
 GSGGIRPKQPRMLQENTRMDMENEESMSYEGDGG EASHGLEDSNISYGSYEEPPDKSNTQDTSFSSIGGYEV
 S EEEEEEEEEQRSGPSVLSQVHLS EDEEDSEDFHSIAGDSLDSDE

Séquence 2 de longueur 652:

>sp|Q9NPI1|BRD7_HUMAN Bromodomain-containing protein 7 OS=Homo sapiens OX=9606 GN=BRD7 PE=1 SV=1

MGKKHKKHKS D KHL YEEYVEKPLKLVLKVGGNEVTELSTGSSGHDSSLFEDKNDHDKHKDRKRKRKKGEKQI
 PGEEKGRKRRRVKEDKKKRDRDRVENEAEKDLQCHAPVRLDLPPEKPLTSSLAKQEEVEQTPLQEALNQLMRQ
 LQRKDPSAFFSFPVTDFIAPGYSMIIKHPMDFSTMEKIKNN DYQSIEELKDNFKLMCTNAM IYKNETIYYK
 AAKKLLHSGMKILS QERIQSLKQSIDFMADLQKTRKQKDGTDTSQSGEDGGCWQREREDSGDAEHAFAKSPSK
 ENKKKDKMDLEDKFKSNNLEREQEQLDRIVKESGGKLT RRLVNSQCEFERRKPDGTTTLGLLHPVDPIVGEPG
 YCPVRLGTTGRQLQSGVNTLQGFKEDKRNVTPVLYLNYG PYSSYAPHYDSTFANISKDDSDLIYSTYGEDSD
 LPSDFS IHEFLATCQDYPYVMADSLDLVLTGGH SRTLQEMEMSLPEDEGHTRTLDTAKEMEITEVEPPGRLD
 SSTQDR LIALKAVTNFGVPVEVFDSEEAEIFQKKLDETTRLLRELQEAQNERLSTRPPNMICLLGPSYREMH
 LAEQVTNNL KELAQQVTPGDIVSTYGVKAMGISIPSPVMENNFVDLTEDTEEPKKT DVAECGPGGS

matrice de substitution utilisée: blosum62.txt

Régularité de gap opening: T = 12, L = 5

Nous constatons que les 2 séquences ayant le taux de similarité le plus élevé, ne sont pas de la même famille :

- La première séquence est un "Transcription initiation factor TFIID subunit 1", et avec un gène de type TAF1
- La deuxième séquence est un "bromodomaine contenant des protéines 7", avec un gène de type BRD7

Peut-être ont-ils des fonctionnalités similaires au sein de l'organisme humain ?

Séquences 1 et 2 du fichier protein-sequences

```
In [36]: 1 L3 = LocalAlignment(3, 12, 2, "blosum62.txt", P2.get_seq(0), P2.get_seq(1)
2 L3.Smith_Waterman()
3
```

Séquence 1 de longueur 598:

```
>sp|Q9H8M2|BRD9_HUMAN Bromodomain-containing protein 9 OS=Homo sapiens OX=9606 GN=BRD9 PE=1 SV=2
```

```
MGKKKKKKHKAERSSYEDYADKPLEKPLKLVLVKVGSGSEVTELSGSGHDSSYYDDRSHERERHKEKKKKKKKK
SEKEKHLDDERRKRKEEKRREREHCDTEGEADDFDPGKKVEVEPPDRPVRACRTQPAENESTPIQQLLE
HFLRQLQRKDPHGFFAFPVTDIAIPGYSMIIKHPMDFGTMKDKIVANEYKSVTEFKADFKLMCDNAMTYNRPD
TVYYKLAKKILHAGFKMMSKQAALLGNEDTAVEEPVPEVVPVQVETAKKSKKPSREVISCMEPEGNACSLTD
STAEHVLALVEHADEARDRINRFLPGGKMGYLRNGDGSLLYSVVNTAEPDADEEETHPVDLSSLSSKLLP
GFTTLGFKDERRNKVTFLSSATTALSMQNNSVFGDLKSDMELLYSAYGDETGVCALSLQEFVKDAGSYSKK
VVDDLQDITGGDHSRTLFLQKQRRNPMKPPDEAKVGDTLGDSSSSVLEFMSMKSYPDVSDISMLSSLGKV
KKELPDDSHLNLDETTLKQLDLHEAQERGGSRPSSNLSSLSNASERDQHHLGSPSRLSVGEQPDVTHDPYE
FLQSPEPAASAKT
```

Séquence 2 de longueur 1363:

```
>sp|060885|BRD4_HUMAN Bromodomain-containing protein 4 OS=Homo sapiens OX=9606 GN=BRD4 PE=1 SV=2
```

```
MSAESGPGTRLRNLPMVDGLETSMSTTQAQAQPANAASTNPPPPETSNNPNPKRQTNQLQYLLRVVLKT
LWKHQFAWPFQQPVDVAVKLNLPDYYKIIKTPMDGTIKRLENNYYWNAQECIQDFNTMFTNCYIYNKPGDDI
VLMAEAELEKFLQKINELPTEETEIMIVQAKGRGRKGTGTAKPGVSTVPNTTQASTPPQTQTPQNPPPVQ
ATPHFPFAVTPDLIVQTPVMTVVPVPPQLQTPPPVPPQPPPPAPAPQPVQSHPPIIAATPQPVKTKKGVKRKA
DTTPTTTIDPIHEPPSLPPEPKTTKLQGRRESSRPVKPKKDVDPDSQQHPAPEKSSKVSEQLKCCSGILKEMF
AKKHAAYAWPFYKPVDEALGLHDYCDIIKHPMDMSTIKSKLEAREYRDAQEFQADVRLMFSNCKYNPPDHE
VVAMARKLQDVFEMRFAPMPDEPEEPVAVSSPAVPPPTKVAPPSSSDSSSDSSSDSSSDSDSDTDDSEEEAQR
AELQEQLKAVHEQLAALSQPQONKPKKKEKDKKEKKKKEKHKRKEEVEENKSKAKEPPPKTKKNNSSNSNVS
KKEPAPMKSKPPPTYESEEEDKCKPMSYEEKRQLSLDINKLPGKELGRVVIHQSRPSLKNPNDEIEIDFE
TLKPSTLRELYVTSCLRKRKPQAEKVDVIAGSSKMGFSSSESESSSESSSDSEDSETEMAPKSKKKGH
PGREQKHHHHHHHQQMQAPAPVPQPPPPPPQPPPPPPPPQPPPPPPPPPPPPPPPPPPPPPPPPPPPPPP
QVPVLEPQLPGSVFDPIGHFTQPIHLHPQELPPHLPQPEHSTPHLNLQHAVVSPPALHNLALPQQPSRPSNR
AAALPPKPARPPAVSPALQTPLLPQPPMAQPPQVLLLEDEPPAPPLTSMQMLYLQQLQKVPPTPLLPVSK
VQSQPPPLPPPPHPSVQQQLQQQPPPPPPPPQPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPPP
PPGQPPPPPPQPAKPPQVIQHHHSRHHKSDPYSTGHLREAPSPLMIHSPQMSQFQSLTHQSPPPQNVQPKKQE
LRAASVVQPQLVVVKEEKIHSPIIRSEPFSPSLRPEPPKHPESIKAPVHLRQPEMKPVDVGRPVIRPPEQN
APPGAPDKDKQKQEPKTPVAPKDKLKIKNMGSWASLVQKHPTTPSSTAKSSSDSFEQFRRAAREKEEREKAL
KAQAEHAEKEKERLRQERMRSREDEDALEQARRAHEEARRRQEQQQQQRQEQQQQQQQAAAVAAAATPQAQS
SQPQSMLDQRELARKREQERRRREMAATIDMNFQSDLLSIFEENLF
```

matrice de substitution utilisée: blosum62.txt

Pénalité de gap affine: I = 12 | E = 2

>>> Alignement:

```
PPDRPVRACRTQPAENESTPIQQLLE---HFLRQLQRKDPHGFFAFP---V-TDAIAPG
:: . : . .:: .:: . . . : . . . . : : . .:: : . . . .
PPKDVDPDSQQHPAPEKSSKVSEQLKCCSGILKEMFAKK-HAAYAWPFYKPVDEALGLH
```

```
-YSMIIKHPMDFGTMKDKIVANEYKSVTEFKADFKLMCDNAMTYNRPD TVYYKLAKKILH
: : : : : : : : : : : : : : : : : : : : : : : : : : : :
DYCDIIKHPMDMSTIKSKLEAREYRDAQEFQADVRLMFSNCKYNPPDHEVVAMARKLQD
```

```
AGFKMMSKQAALLGNEDTAVEEPVPEVVPVQVETAKKSKKPSREVISCMEPEGNACSLT
. .:: . . : . : : : : : : : : : : : : : : : : : : : : :
V-FEM--RFAKMP--DEP-EEPVVAVSSPAVPPPTKVAPPSSSDSSS-DSSSDSDSST
```

```
DSTAEHVLALVE
... .:: .::
DDSEEEAQRALAE
```

==> similarité: 58.03108808290155 %

==> Identité: 31.606217616580313 %

```

Algorithm: Smith-Waterman (SSE2, Michael Farrar 2006) (7.2 Nov 2010)
Parameters: data/blosum62.mat matrix (11:-4), open/ext: -12/-2
Scan time: 0.150

>>unknown 1362 bp (1362 aa)
Waterman-Eggert score: 144; 52.2 bits; E(1) < 1.6e-10
35.6% identity (64.4% similar) in 118 aa overlap (121-229:330-446)

      130      140      150      160      170
unknown PPDRPVACRTQPAENESTPIQQLLE---HFLRQLQRKDPHGFFAFPV---TDAIAPG--
      ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::
unknown PPKKDVPDSQQHPAPEKSSKVSEQLKCCSGILKEMFAKK-HAAYAWPFYKPVDEALGLH
      330      340      350      360      370      380

      180      190      200      210      220
unknown -YSMIIKHPMDFGTMKDKIVANEYKSVTEFKADFKLMCDNAMTYNRPDTVYYKLAKKI
      :  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::
unknown DYCDIHKHPMDSTIKSLEAREYRDAQEFQADVRLMFSNCKYKYNPPDHEVVMARKL
      390      400      410      420      430      440

>--
Waterman-Eggert score: 108; 40.4 bits; E(1) < 5.8e-07
33.3% identity (60.4% similar) in 111 aa overlap (111-218:36-142)

      120      130      140      150      160      170
unknown DPGKKVEVEPPDRPVACRTQPAENESTPIQQLLEHFLRQLQRKDPHGFFAFPVTDIAIA
      ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::
unknown QPANAASTNPPP--PETSNNPKP-KRQTNQLQYLLRVVLKTLWKHQFAWPFQQPV-DAVK
      40      50      60      70      80      90

      180      190      200      210
unknown ---PGYSMIIKHPMDFGTMKDKIVANEYKSVTEFKADFKLMCDNAMTYNRP
      :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :
unknown LNLPDYYKIIKTPMDMGTIKKRLNNYYWNAQECIQDFNTMTNICYIYNKP
      100      110      120      130      140

>--
Waterman-Eggert score: 84; 32.5 bits; E(1) < 0.00014
26.5% identity (60.2% similar) in 98 aa overlap (27-124:476-572)

      30      40      50      60      70      80
unknown PLKLVLVKVGSEVTELSGSGHDSSYYDDRSDBHERERHKEKKKKKKKSEKEHLDDEERR
      :  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::
unknown PTKVVAPPSSSDSSSDSSSDSSST-DDSEEEAQLAELQEQLKAVHEQLAALSQPPQN
      480      490      500      510      520      530

      90      100      110      120
unknown KRKEEKKRKREREHCDTEGEADDFDPGKKVEVEPPDR
      :  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::  ::
unknown KPKKKEKDKKEKKKEKHKEEVEENKSKAKEPPPKK
      540      550      560      570

```

Nous constatons que les résultats ne sont pas les mêmes que ceux de LALIGN. Les 2 résultats ont par ailleurs utilisées les mêmes matrices de substitution BLOSUM62 et des pénalités de gap affine $I = 12$ et $E = 2$

Notre alignement contient le sous alignement que nous donnons LALIGN. Par exemple, Notre première solution part de "PPDR" et finit en "ALVE" tandis que la première solution de LALIGN s'arrête bien avant en "AKKI". La position des gap change également. Par exemple, après la section 160 de la première solution, nous constatons que le V est suivi de 3 gaps dans LALIGN tandis que notre V est précédé de plusieurs gap.

La 3ème solution par contre, est identique à celle de LALIGN à peu près. Dans celui-ci, aucun gap n'a été introduit du côté de LALIGN et 2 sont apparus dans notre solution.

D'après les données d'UNIPROT, ces 2 séquences sont des bromodomains mais contenant des protéines de type différentes (9 pour la première et 4 pour la deuxième. Ces séquences sont donc relativement similaires.

Conclusion

Pour conclure, notre problème était d'aligner des séquences d'acides aminés de sorte à trouver leurs similitudes maximales dissimulées et déterminer si celles-ci sont homologues (2 gènes ou 2 protéines qui partagent un ancêtre commun). Pour ce faire, nous avons utilisé les algorithmes de Needleman-Wunsch et Smith-Waterman pour les alignements globaux et locaux.

Nous avons constaté que fatalement, lorsque 2 séquences appartiennent à la même famille (2 bromodomaines contenant des protéines 4 par exemple) ainsi que la même gène, l'alignement global ou local donnait un taux de similarité élevé. D'autre part, certains alignements concernaient 2 séquences n'étant pas de la même famille. Le taux de similarités était tout de même élevé, ce qui laisse croire qu'ils ont des fonctionnalités similaires.