

2I003 Initiation à l'Algorithmique
Projet: Alignement des séquences ADN

Andrea KOSTAKIS

Novembre 2014

Première partie

Distance minimum de l'alignement de deux séquences

1 Exercice 1

Rappelons les différentes définitions utiles pour la suite du projet :

Une séquence d'ADN (acide Deoxyribonucleique) est une suite de nucléotides. Chaque nucléotide est désigné par la première lettre de sa molécule de base.

Une séquence d'ADN est ainsi un mot sur l'alphabet $\mathcal{A} = \{A, C, G, T\}$. Soient deux séquences ADN composées respectivement de n, m nucléotides : $a = a_0 \dots a_{n-1}$ et $b = b_0 \dots b_{m-1}$, alors

$\forall i \in \{0, \dots, n-1\}, a_i \in \mathcal{A}$ et $\forall j \in \{0, \dots, m-1\}, b_j \in \mathcal{A}$.

Alignement : On appelle alignement de longueur $L \geq \max(n, m)$ des séquences a et b tout couple (a^*, b^*) où $a^* = a_0^* \dots a_L^*$ et $b^* = b_0^* \dots b_L^*$ sont deux mots de longueur L construits en insérant des éléments nuls dans les séquences a et b respectivement. C'est à dire, $\forall i \in \{0, \dots, L-1\}$, a_i^* et b_i^* sont des éléments de $\mathcal{A} \cup \{-\}$.

Distance de Levenshtein : la distance entre deux éléments de $\mathcal{A} \cup \{-\}$ est définie par :

- $D(\alpha, \alpha) = 0, \forall \alpha \in \mathcal{A} \cup \{-\}$
- $D(\alpha, \beta) = 1, \forall (\alpha, \beta) \in (\mathcal{A} \cup \{-\})^2$ avec $\alpha \neq \beta$

La distance d'un alignement est alors :

$$D(a^*, b^*) = \sum_{k=0}^{L-1} D(a_k^*, b_k^*)$$

1.1 Alignements : Distance et longueur

Montrons que, si pour $j \in \{0, \dots, L-1\}, a_j^* = b_j^* = -$, on peut construire un autre alignement (a'^*, b'^*) de même distance et longueur $L-1$.

Soit L la longueur de l'alignement (a^*, b^*) si pour $j \in \{0, \dots, L-1\}$, $a_j^* = b_j^* = -$, alors la distance $D(a^*, b^*)$ est :

$$D(a^*, b^*) = \sum_{k=0}^{L-1} D(a_k^*, b_k^*) = \sum_{k=0}^{j-1} D(a_k^*, b_k^*) + D(a_j^*, b_j^*) + \sum_{k=j+1}^{L-1} D(a_k^*, b_k^*)$$

par définition, $D(a_j^*, b_j^*) = 0$ donc,

$$D(a^*, b^*) = \sum_{k=0}^{j-1} D(a_k^*, b_k^*) + \sum_{k=j+1}^{L-1} D(a_k^*, b_k^*)$$

Construisons un autre alignement (a'^*, b'^*) de longueur L' en retirant les éléments a_j^* et b_j^* (càd l'élément nul - en position j de a^* et b^*) de l'alignement (a^*, b^*) , on obtient donc l'alignement (a'^*, b'^*) où

$$a'^* = a_0^* \dots a_{j-1}^* a_{j+1}^* \dots a_{L-1}^*$$

$$b'^* = b_0^* \dots b_{j-1}^* b_{j+1}^* \dots b_{L-1}^*$$

On remarque d'après la construction que l'on obtient un alignement de longueur $L' = L - 1$. La distance de (a'^*, b'^*) est :

$$D(a'^*, b'^*) = \sum_{k=0}^{L'-1} D(a_k'^*, b_k'^*)$$

par constuction :

$$= \sum_{k=0}^{j-1} D(a_k^*, b_k^*) + \sum_{k=j+1}^{L-1} D(a_k^*, b_k^*)$$

$$= D(a^*, b^*)$$

On obtient donc bien un autre alignement (a'^*, b'^*) de même distance et de longueur $L - 1$.

1.2 Alignement optimal de sous séquence

Soit le couple (a^*, b^*) de longueur L des séquences $a = a_0 \dots a_{n-1}$ et $b = b_0 \dots b_{m-1}$ de distance minimum avec $n > 0$ et $m > 0$. $c(n-1, n) \geq 1 + c(n-2, n-1) + c(n-2, n) + c(n-1, n-1)$

1.2.1 Demonstration

Montrons par l'absurde que si $a_{L-1}^* = -$ et $b_{L-1}^* = b_{m-1}$ alors le couple $(a_0^* \dots a_{L-2}^*, b_0^* \dots b_{L-2}^*)$ est un alignement optimal de a avec la sous séquence et $b_0 \dots b_{m-2}$.

Raisonnons donc par l'absurde :

Supposons que le couple $(a_0^* \dots a_{L-2}^*, b_0^* \dots b_{L-2}^*)$ n'est pas un alignement optimal.

Donc, il existe un autre alignement (a'^*, b'^*) provenant des mêmes séquences tel que

$$D(a'^*, b'^*) < D(a_0^* \dots a_{L-2}^*, b_0^* \dots b_{L-2}^*)$$

Par définition de la distance,

$$D(a'^*, b'^*) < \sum_{k=0}^{L-2} D(a_k^*, b_k^*)$$

(a^*, b^*) est un alignement optimal et sa distance peut s'écrire de la manière suivante :

$$D(a^*, b^*) = \sum_{k=0}^{L-2} D(a_k^*, b_k^*) + D(a_{L-1}^*, b_{L-1}^*)$$

on remarque que

$$D(a_{L-1}^*, b_{L-1}^*) = 1 \text{ car } (a_{L-1}^* \neq b_{L-1}^*)$$

enfin, on a également

$$\sum_{k=0}^{L-2} D(a_k^*, b_k^*) = D(a_0^* \dots a_{L-2}^*, b_0^* \dots b_{L-2}^*)$$

on peut donc écrire

$$D(a_0^* \dots a_{L-2}^*, b_0^* \dots b_{L-2}^*) = D(a^*, b^*) - 1$$

A fortiori,

$$D(a'^*, b'^*) < D(a^*, b^*) - 1$$

Cette inégalité nous ramène à une contradiction. (en retirant une unité à une distance optimale on obtient forcément un résultat optimal). Il n'existe donc pas d'autre alignement optimal. On peut en conclure que le couple $(a_0^* \dots a_{L-2}^*, b_0^* \dots b_{L-2}^*)$ est un alignement optimal de a avec la sous séquence et $b_0 \dots b_{m-2}$.

1.2.2 Autres valeurs possibles

Examinons le cas des autres valeurs possible de a_{L-1}^* et b_{L-1}^* :

- si $a_{L-1}^* = -$ et $b_{L-1}^* = -$ alors, $D(a_0^* \dots a_{L-2}^*, b_0^* \dots b_{L-2}^*) = D(a^*, b^*)$ (question 1.1) et $(a_0^* \dots a_{L-2}^*, b_0^* \dots b_{L-2}^*)$ est un alignement optimal
- respectivement si $a_{L-1}^* = b_{L-1}^*$
- si $a_{L-1}^* = a_{n-1}$ et $b_{L-1}^* = -$ alors $(a_0^* \dots a_{L-2}^*, b_0^* \dots b_{L-2}^*)$ est un alignement optimal avec le même raisonnement que la question précédente

1.3 Décompositon de la distance minimale

Pour tout couple $(i, j) \in \{-1, \dots, n-1\} \times \{-1, \dots, m-1\}$ on désigne par $\bar{D}_{i,j}$ la distance minimale d'un alignement des séquences $a_0^* \dots a_i^*$ et $b_0^* \dots b_j^*$. Le cas $i = -1$ (resp. $j = -1$) correspond à une séquence vide notée ϵ . Par convention, on pose

$$D_{-1,-1} = 0$$

Montrons que,

$$\forall i \in \{0, \dots, n-1\}, \bar{D}_{i,-1} = \bar{D}_{i-1,-1} + D(a_i, -)$$

et que,

$$\forall j \in \{0, \dots, m-1\}, \bar{D}_{-1,j} = \bar{D}_{-1,j-1} + D(-, b_j)$$

$\forall i \in \{0, \dots, n-1\}$, $\bar{D}_{i,-1}$ correspond à la distance optimale entre $a = a_0 \dots a_i$ et ϵ . D'après 1.2, si $D_{i,-1}$ est optimale alors $D_{i-1,-1}$ est optimale. on peut donc écrire,

$$\bar{D}_{i,-1} = \bar{D}_{i-1,-1} + D(a_i, -)$$

car $\forall i \in \{0, \dots, n-1\} D(a_i, -)$ est la distance du dernier couple.

Respectivement,

$$\forall j \in \{0, \dots, m-1\}, \bar{D}_{-1,j} = \bar{D}_{-1,j-1} + D(-, b_j)$$

1.4 Alignement optimal : calcul

Démontrons que pour tout couple $(i, j) \in \{0, \dots, n-1\} \times \{0, \dots, m-1\}$,

$$\bar{D}_{i,j} = \min(\bar{D}_{i-1,j} + D(a_i, -), \bar{D}_{i-1,j-1} + D(a_i, b_j), \bar{D}_{i,j-1} + D(b_j, -))$$

Examinons le cas d'un alignement optimal pour tout couple (i, j) , on observe trois cas particuliers que l'on peut décomposer d'après les questions précédentes :

- si $n > m$ alors, $\forall (i, j) \in \{0, \dots, n-1\} \times \{0, \dots, m-1\}$

$$\bar{D}_{i,j} = \bar{D}(a_{i-1}, b_j) + D(a_i, -) = \bar{D}_{i-1,j} + D(a_i, -) \quad (1)$$

- si $n = m$, alors $\forall (i, j) \in \{0, \dots, n-1\} \times \{0, \dots, m-1\}$

$$\bar{D}_{i,j} = \bar{D}_{i-1,j-1} + D(a_i, b_j) \quad (2)$$

- si $n < m$ alors, $\forall (i, j) \in \{0, \dots, n-1\} \times \{0, \dots, m-1\}$

$$\bar{D}_{i,j} = \bar{D}_{i,j-1} + D(b_j, -) \quad (3)$$

Par conséquent si l'on veut obtenir l'alignement optimal, on doit regrouper les trois cas, c'est-à-dire minimiser (1)(2)(3).

D'où,

$$\bar{D}_{i,j} = \min(\bar{D}_{i-1,j} + D(a_i, -), \bar{D}_{i-1,j-1} + D(a_i, b_j), \bar{D}_{i,j-1} + D(b_j, -))$$

C'est une définition récursive de la distance d'un l'alignement optimal.

1.5 DistanceMinRec(a,b)

D'après 1.4 on peut définir la fonction récursive suivante :

Supposons définie la fonction Distance(a[i],b[i]) qui rend la distance entre deux éléments. (implémenté Exercice 2).

```

1 def DistanceMinRec(a,b):
2   n=len(a)
3   m=len(b)
4   d=0
5   if (n==0 and m==0):
6     return 0
7   if (n==0):
8     return m
9   if (m==0):
10    return n
11  return min(DistanceMinRec(a[n-1],b)+1,
12             DistanceMinRec(a[n-1],b[m-1])+ Distance(a[n],b[m]), \ subsection{Fonctions}
13             DistanceMinRec(a,b[m-1])+1)
```

1.6 Terminaison et validité

Par récurrence (forte) sur $k = n + m$ Prouvons que DistanceMinRec se termine et renvoie la distance minimale de l'alignement des séquences a,b.

Base : Pour $n=m=0$ ($k=0$) la fonction se termine et renvoie 0, qui est bien la distance minimale des séquences nulles.

Induction : Soit DistanceMinRec se termine à un certain rang k et renvoie la distance minimale des k derniers éléments de la séquence a, b .

A chaque appel récursif de DistanceMinRec on renvoie la distance des sous-séquences $k-1$ à $n+m$.

Au rang $k_0 + 1$ ($k_0 > k$) la fonction renvoie la distance minimale de la sous-séquence $k-1$ à $n+m$.

A ce rang on obtiendra $m=0$ où $n=0$ donc la fonction se terminera. De plus, comme pour chaque sous-séquence la distance minimale est calculée pour les $k_0 - 1$ éléments.

Conclusion : Par récurrence au rang k_0 DistanceMinRec(a, b) se termine et est valide.

1.7 Complexité de DistanceMinRec

Soit $c(n, m)$ la complexité pour deux séquences de taille respective n et m . On note $u_n = c(n, m)$ pour deux séquences de même taille.

1.7.1 Calcul de $c(n, m)$

Pour $n > 0$ et $m > 0$, d'après la fonction on appelle récursivement DistanceMinRec($a[n-1], b$), DistanceMinRec($a[n-1], b[m-1]$) et DistanceMinRec($a, b[m-1]$) puis Distance($a[n], b[m]$) qui s'exécute en temps constant ($= 1$).

La complexité $c(n, m)$ s'écrit :

$$c(n, m) \geq 1 + c(n-1, m) + c(n-1, m-1) + c(n, m-1)$$

En remplaçant dans l'inégalité obtenue, on obtient :

$$c(n-1, n) \geq 1 + c(n-2, n-1) + c(n-2, n) + c(n-1, n-1)$$

$$\text{et } c(n, n-1) \geq 1 + c(n-1, n-2) + c(n-1, n-1) + c(n-1, n-2)$$

On en déduit que pour $n > 1$:

$$c(n-1, n) \geq c(n-1, n-1) \text{ et } c(n, n-1) \geq c(n-1, n-1)$$

1.7.2 Complexité u_n

$$u_n = c(n, n) \geq 1 + c(n-1, n-1) + c(n-1, n) + c(n, n-1)$$

d'après 1.7.1 :

$$c(n-1, n) \geq c(n-1, n-1)$$

$$c(n, n-1) \geq c(n-1, n-1)$$

donc,

$$u_n \geq 1 + 3c(n-1, n-1)$$

$$u_n \geq 1 + 3u_{n-1} \text{ et } u_0 = 1$$

Donc, comme u_{n-1} est borné par $(-1/3)^{n-1}$ alors $u_n \geq (-1/3)^n$ qui est bien une fonction exponentielle.

1.8 DistanceMinIter(a,b)

Voici le code proposé :

```
1 DistanceMonIter(a,b):
2 n=len(a)
3 m=len(b)
4 result=0
5 for (i=0;i<=n;i++):
6   M[i][0]=i
7 for (j=0;j<=m;j++):
8   M[0][j]=j
9 for (i=0;i<=n;i++):
10  for (j=0;j<=m;j++):
11    flag=Distance(a[n],b[m])
12    M[i][j]=min(M[i-1][j]+1,M[i-1][j-1]+val,M[i][j-1])
13 result=M[n][m]
14 return result
```

1.9 Complexité de DistanceMinIter(a,b)

La fonction itérative ci-dessus est composée de plusieurs boucles de types for. On parcourt $n \times m$ cases car M est un tableau à deux dimensions.

D'où une complexité en $\Theta(n \times m)$

Deuxième partie

Mesure expérimentale de la complexité

2 Exercice 2

Le langage choisit pour la suite du projet est le C. Rappelons que la complexité expérimentale des fonctions est la mesure du temps d'exécution.

2.1 Fonctions

Nous créons les fonctions Distance(a,b) et minimum(a,b,c) supplémentaires afin d'alléger le code.

Les listings de codes se trouvent à la fin du rapport. Voici un exemple d'exécution pour des séquences de longueur $L = 10$ pour la fonction itérative et récursive respectivement le retour est en première ligne le temps d'exécution puis le résultat :

```
1 La distance minimale des sequences
2 CCTCTGTCCC
3 et
4 TGCTCTATGT
5 = 6
6 Rec= 0.970000sec
7 La distance minimale des sequences
8 CCTCTGTCCC
9 et
10 TGCTCTATGT
11 = 6
```

Nous créons également les fichiers resultatsiter.txt et resultasrec.txt avec le temps d'exécution (1ère colonne) et longueur (2ème colonne) :

```
1 ==> resultatsite.txt <==
2 0.000183 10
3 0.000268 15
4 0.000090 20
5 0.000572 40
6 0.000566 50
7 0.001793 100
8 0.004398 200
9 0.009176 300
10 0.015309 400
11 0.023304 500
```

```

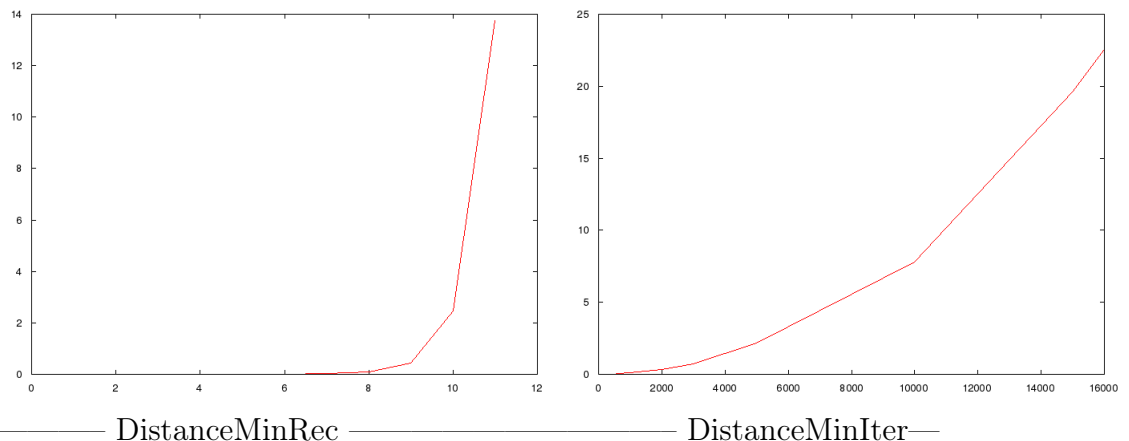
12 0.032963 600
13 0.084551 1000
14 0.321449 2000
15 0.714503 3000
16 2.187023 5000
17 7.810801 10000
18 19.658818 15000
19 22.574136 16000
20 ==> resultatsrec.txt <==
21 0.000008 1
22 0.000019 2
23 0.000053 3
24 0.000200 4
25 0.000589 5
26 0.003614 6
27 0.016190 7
28 0.080566 8
29 0.438319 9
30 2.449000 10
31 13.736783 11

```

2.2 Temps d'exécution

En fonction des performance de l'ordinateur on observe que pour des longueurs supérieures à 12 la pile de récursion ne marche pas, le processus est tué par le système. Pour la fonction itérative cette limite de mémoire est un peu près à L=20000.

D'après les fichiers de résultat et la fonction gnuplot du shell on obtient les graphes suivants en fonction de la taille.



D'après les graphes on voit bien que DistanceMinRec à une complexité expérimentale exponentielle. DistanceMinIter à une complexité expérimentale linéaire si l'on néglige les "bruits du processeur".

2.3 Listings

```
1  ==> Distance.c <==
2  #include <stdio.h>
3  #include <stdlib.h>
4  #include <string.h>
5  #include <time.h>
6
7  int Distance(char a, char b)
8  {
9      if (a==b)
10         return 0;
11     return 1;
12 }
13
14
15 ==> DistanceMinIter.c <==
16 #include <stdio.h>
17 #include <stdlib.h>
18 #include <string.h>
19 #include <time.h>
20
21 #include "Distance.h"
22 #include "Minimum.h"
23
24 int DistanceMinIter(char* a, char* b)
25 {
26
27     int i, j;
28     int n=strlen(a);
29     int m=strlen(b);
30
31     int* M=(int *) malloc((n+1)*(m+1)*sizeof(int));
32     int tmp;
33     int res=0;
34
35     for(i=0; i <=n ; i++){
36         *(M + i*(m+1))=i;
37     }
38     for(j=0; j <=m ; j++){
39         *(M+j)=j;
40     }
41
42
43     for(i=1; i <=n ; i++){
44         for(j=1 ; j<=m ; j++){
45             tmp=Distance(*(a+i-1),*(b+j-1));
46             *(M + i*(m+1) + j) = minimum(*(M +(i-1)*(m+1) + j) +1,
47                                     *(M +(i-1)*(m+1) + (j-1))+tmp,
```

```

48             *(M + i*(m+1) + (j-1)) + 1);
49     }
50     printf("\n");
51 }
52 res = *(M + n*(m+1) + m);
53
54 free(M);
55 return res;
56
57 }
58
59 ==> DistanceMinRec.c <==
60 #include <stdio.h>
61 #include <stdlib.h>
62 #include <string.h>
63 #include <time.h>
64
65 #include "Distance.h"
66 #include "Minimum.h"
67
68
69 int DistanceMinRec(char* a, char* b)
70 {
71
72     int n = strlen(a);
73     int m = strlen(b);
74
75     if (n == 0 && m == 0)
76         return 0;
77     if (n == 0)
78         return m;
79     if (m == 0)
80         return n;
81
82     char* a_bis = strdup(a, n-1); //strdup: duplicate a string,
83                                     rend un pointeur sur le nouveau
84     char* b_bis = strdup(b, m-1);
85
86     return minimum(DistanceMinRec(a_bis, b) + 1,
87                   DistanceMinRec(a_bis, b_bis) +
88                   Distance(*(a+n-1), *(b+m-1)),
89                   DistanceMinRec(a, b_bis) + 1);
90
91 }
92
93
94 ==> main.c <==
95
96 #include <stdio.h>

```

```

97 #include <stdlib.h>
98 #include <string.h>
99 #include <time.h>
100
101 #include "Distance.h"
102 #include "Minimum.h"
103 #include "SequenceAleatoire.h"
104 #include "DistanceMinRec.h"
105 #include "DistanceMinIter.h"
106
107 #define L 10
108
109 int main(){
110     srand(time(NULL));
111     char* a=SeqAleatoire(L);
112     char* b=SeqAleatoire(L);
113
114     int dm;
115     int dmin;
116
117     clock_t start,end;
118
119
120     FILE *fit , *frec;
121     fit=fopen("resultatsite.txt","a+");
122     frec=fopen("resultatsrec.txt","a+");
123     start = clock();
124     {
125         dm=DistanceMinIter(a,b);
126     }
127     end = clock();
128     printf("Iter=%fsec\n", (end - start)/((double)CLOCKS_PER_SEC);
129     fprintf(fit, "%f%d\n", (end - start)/((double)CLOCKS_PER_SEC, L);
130     printf("La distance minimale des s[U+FFFD]mes\n%s\nnet\n%s\n=%d\n",
131         a,b,dm);
132
133     start = clock();
134     {
135         dmin=DistanceMinRec(a,b);
136     }
137     end = clock();
138
139     printf("_Rec=%fsec\n", (end - start)/((double)CLOCKS_PER_SEC);
140     fprintf(frec, "%f%d\n", (end - start)/((double)CLOCKS_PER_SEC, L);
141     printf("La distance _minimale des s[U+FFFD]mes\n%s\nnet\n%s\n=%d\n",
142         a,b,dmin);
143
144
145     free(a);

```

```

146     free(b);
147
148     return 0;
149 }
150
151 ==> Minimum.c <==
152 #include <stdio.h>
153 #include <stdlib.h>
154 #include <string.h>
155 #include <time.h>
156
157
158 int minimum(int a,int b,int c)
159 {
160     if(a==b && a==c)
161         return a;
162     if(a<=b && a<=c)
163         return a;
164     if(b<=a && b<=c)
165         return b;
166
167     return c;
168 }
169
170
171 ==> SequenceAleatoire.c <==
172 #include <stdio.h>
173 #include <stdlib.h>
174 #include <string.h>
175 #include <time.h>
176
177
178
179 char* SeqAleatoire(int n)
180 {
181     int flag;
182     int i;
183     char* chaine=malloc((n+1)*sizeof(char));
184
185
186     for(i=0;i<n;i++)
187     {
188         flag=rand()%4;
189         if(flag==0)
190             *(chaine+i)='A';
191         if(flag==1)
192             *(chaine+i)='C';
193         if(flag==2)
194             *(chaine+i)='G';

```

```
195     if (flag==3)
196         *(chaine+i)='T';
197     }
198
199     *(chaine+i)='\0';
200     return chaine;}
```