

Intérêt des séquences

- La séquence nucléotidique d'un gène détermine la séquence d'acides aminés de la protéine.
- La séquence d'une protéine détermine sa structure et sa fonction
- Généralement, une similarité de séquence implique une similarité de structure et de fonction (l'inverse n'est pas toujours vrai).
- Évolution basée, en grande partie, sur la duplication suivie de modification
⇒ beaucoup de redondance dans les banques de données

Problématiques

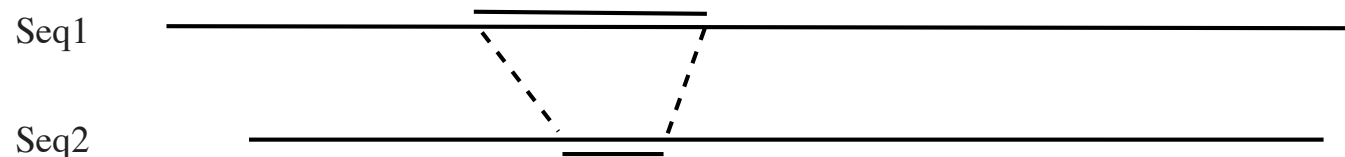
- Est-ce qu'une nouvelle séquence a déjà été complètement ou partiellement déposée dans une banque de données?
- Est-ce que cette séquence contient un gène?
- Est-ce que ce gène appartient à une famille connue?
- Existe-t-il d'autres gènes homologues?
- Comment ce gène a-t-il évolué par rapport aux autres gènes homologues déjà identifiés?

Alignement de séquences

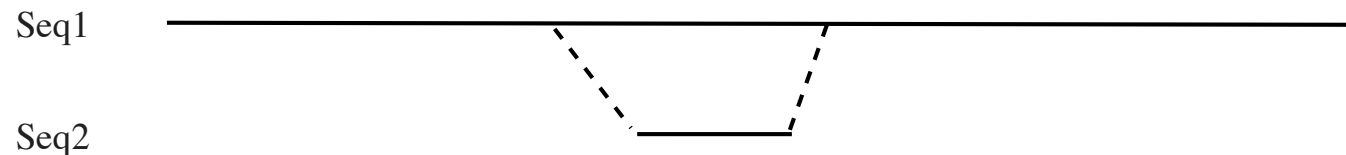
Alignement global: Deux séquences de protéines appartenant à la même famille, études phylogénétiques



Alignement local: Deux séquences de protéines appartenant à des familles différentes mais ayant un ou des domaines communs



Recherche de motif: Recherche des occurrences d'un gène dans un génome



Distance d'édition

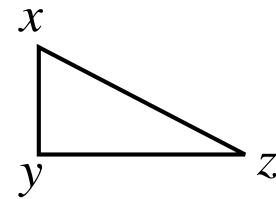
Pour comparer des séquences, on va définir une distance:

Définition: Une distance D est une relation ayant les propriétés suivantes:

$$D(x, x) = 0$$

$$D(x, y) = D(y, x)$$

$$D(x, z) \leq D(x, y) + D(y, z) \quad \text{inégalité du triangle}$$



Distance naturelle: compter le nombre d'**insertions**, de **suppressions** et de **substitutions** nécessaire pour passer d'une séquence à une autre.

Distance d'édition (suite)

Exemple: $S_1 = CATAGTG$ $S_2 = GTCAGGT$

S	Su	M	I	M	I	M	M	Su
C	A	T		A		G	T	G
G		T	C	A	G	G	T	

Distance d'édition entre S_1 et S_2 : Nombre minimal d'insertions, suppressions et substitutions nécessaire pour transformer S_1 en S_2

Une insertion/suppression est représentée par un '-':

C	A	T	-	A	-	G	T	G
G	-	T	C	A	G	G	T	-

Alignement global

Il existe plusieurs alignements possibles étant donné 2 séquences.

Comment trouver un alignement **optimal** i.e un alignement ayant une distance d'édition **minimal**??

Idée naïve: Énumérer tous les alignements possibles pour les 2 séquences et en choisir un dont la distance d'édition est minimal.

Exemple: $S_1 = AC$ et $S_2 = AGC$

Les alignements possibles ici sont:

$$\left\{ \begin{array}{cccccccc} AC- & A-C & -AC & AC-- & AC-- & A--C & A-C- & -A-C \\ AGC' & AGC' & AGC' & -AGC' & A-GC' & -AGC' & AG-C' & AGC-', \dots \end{array} \right\}$$

Combien d'alignements?

Si $f(n,m)$ est le nombre d'alignements entre une séquence $a = a_1 \dots a_n$ de n lettres et une autre séquence $b = b_1 \dots b_m$ de m lettres alors, on a:

$$f(0,0) = 1$$

$$f(0,m) = 1$$

$$f(n,0) = 1$$

$$f(n,m) = \begin{array}{l} \text{Nombre d'alignements de } a_1 \dots a_{n-1} \quad a_n \\ \text{avec } b_1 \dots b_m \quad \text{---} \\ + \\ \text{Nombre d'alignements de } a_1 \dots a_{n-1} \quad a_n \\ \text{avec } b_1 \dots b_{m-1} \quad b_m \\ + \\ \text{Nombre d'alignements de } a_1 \dots a_n \quad \text{---} \\ \text{avec } b_1 \dots b_{m-1} \quad b_m \end{array}$$

$$\Rightarrow f(n,m) = f(n-1,m) + f(n-1,m-1) + f(n,m-1)$$

Combien d'alignement?

$$\Rightarrow f(n,m) = f(n-1,m) + f(n-1,m-1) + f(n,m-1)$$

	0	1	2	3	4	5
0	1	1	1	1	1	1
1	1	3	5	7	9	11
2	1	5	13	25	41	61
3	1	7	25	63	129	231
4	1	9	41	129	321	681

En fait, on peut montrer que

$$f(n,n) \sim (1 + \sqrt{2})^{2n+1} \cdot \sqrt{n}$$

$$\Rightarrow f(1000, 1000) \sim 10^{764}$$

Programmation dynamique

Étant donné 2 séquences, $S = s_1s_2 \dots s_m$ et $T = t_1t_2 \dots t_n$, on définit

$D(i,j)$: distance d'édition entre le préfixe de taille i de S , $s_1 \dots s_i$, et le préfixe de taille j de T , $t_1 \dots t_j$.

D définit une matrice de taille $(m+1) \times (n+1)$ qu'on appelle la matrice de programmation dynamique

L'idée est alors d'exprimer $D(i,j)$ en fonction des valeurs de D pour des paires d'indices plus petits que (i,j)

Calculer $D(i,j)$ à partir des 3 cases $(i-1,j)$, $(i, j-1)$ et $(i-1,j-1)$:

1. L'alignement se termine par la suppression de s_i

$$\begin{array}{ccc} \text{Alignement optimal de } s_1 \dots s_{i-1} & & s_i \\ \text{avec } t_1 \dots t_j & & - \\ D(i-1,j) & + & 1 \end{array}$$

2. L'alignement se termine par l'insertion de t_j

$$\begin{array}{ccc} \text{Alignement optimal de } s_1 \dots s_i & & - \\ \text{avec } t_1 \dots t_{j-1} & & t_j \\ D(i,j-1) & + & 1 \end{array}$$

3. L'alignement se termine par l'alignement de s_i avec t_j

$$\begin{array}{ccc} \text{Alignement optimal de } s_1 \dots s_{i-1} & & s_i \\ \text{avec } t_1 \dots t_{j-1} & & t_j \\ D(i-1,j-1) & + & \begin{array}{l} 1 \text{ si } s_i \neq t_j \\ 0 \text{ sinon} \end{array} \end{array}$$

Remplissage de la table

Conditions initiales: $D(i, 0) = i, \quad \forall i \quad 0 \leq i \leq m$
 $D(0, j) = j, \quad \forall j \quad 0 \leq j \leq n$

Relation de récurrence pour $i, j > 0$:

$$D(i, j) = \min \begin{cases} D(i-1, j) & +1 \\ D(i, j-1) & +1 \\ D(i-1, j-1) + \delta(i, j) \end{cases},$$

Où $\delta(i, j) = 0$ si $x_i = y_j$ et 1 sinon.

Complexité: Pour remplir chaque case de la table, on examine 3 cases.
Il y a $O(nm)$ cases et donc complexité en temps de $O(nm)$

Exemple

Relation de réurrence pour $i, j > 0$:

Conditions initiales: $D(i, 0) = i, \quad \forall i \quad 0 \leq i \leq m$

$D(0, j) = j, \quad \forall j \quad 0 \leq j \leq n$

$$D(i, j) = \min \begin{cases} D(i-1, j) & +1 \\ D(i, j-1) & +1 \\ D(i-1, j-1) + \delta(i, j) \end{cases}$$

Où $\delta(i, j) = 0$ si $x_i = y_j$ et 1 sinon.

Exemple: Calculez la valeur d'un alignement global entre les mots ACGTTA et AACTA

	-	A	C	G	T	T	A
-	0	1	2	3	4	5	6
A	1	0	1	2	3	4	5
A	2	1	1	2	3	4	4
C	3	2	1	2	3	4	5
T	4	3	2	2	2	3	4
A	5	4	3	3	3	3	3

Trouver un alignement optimal

Au cours du remplissage de la table, garder des **pointeurs**:

- de (i,j) à $(i-1,j)$ si $D(i,j) = D(i-1,j) + 1$
- de (i,j) à $(i,j-1)$ si $D(i,j) = D(i,j-1) + 1$
- de (i,j) à $(i-1,j-1)$ si $D(i,j) = D(i-1,j-1) + \delta(i,j)$

Un alignement optimal: Commencer à la case (m,n) et suivre des pointeurs jusqu'à la case $(0,0)$.

Une case peut contenir plusieurs pointeurs: **plusieurs alignements optimaux** possibles

Exemple: Calculez la valeur d'un alignement global entre les mots ACGTTA et AACTA

	-	A	C	G	T	T	A
-	0	1	2	3	4	5	6
A	1	0	1	2	3	4	5
A	2	1	1	2	3	4	4
C	3	2	1	2	3	4	5
T	4	3	2	2	2	3	4
A	5	4	3	3	3	3	3

- A C G T T A
A A C - - T A

Distance versus similarité

Plutôt que de mesurer la différence entre 2 séquences, mesurer leur degré de similarité

$P(a,b)$: score de l'appariement (a,b) . Positif si $a=b$ et ≤ 0 sinon

$V(i,j)$: valeur de l'alignement optimal entre $s_1 \dots s_i$ et $t_1 \dots t_j$

Conditions initiales: $V(i,0) = \sum_{1 \leq k \leq i} P(s_k, -)$, $V(0,j) = \sum_{1 \leq k \leq j} P(-, t_k)$,

Relation de récurrence:

$$V(i,j) = \max \begin{cases} V(i, j-1) & + P(-, t_j) \\ V(i-1, j) & + P(s_i, -) \\ V(i-1, j-1) & + P(s_i, t_j) \end{cases}$$

Algorithme de Needleman-Wunch (1970)

Si le score pour une insertion ou suppression est de -2
 le score pour une substitution est de -1 et
 le score pour un match est de 2 :

Exemple: Calculez la similarité maximale entre les mots ACGTTA
 et AACTA

	-	A	C	G	T	T	A
-	0	-2	-4	-6	-8	-10	-12
A	-2	2	0	-2	-4	-6	-8
A	-4	0	1	-1	-3	-5	-4
C	-6	-2	2	0	-2	-4	-6
T	-8	-4	0	1	2	0	-2
A	-10	-6	-2	-1	0	1	2

Recherche approché d'un motif

Problème: On a un “petit” motif P de taille m et une “longue” séquence T de taille n et on veut trouver toutes les occurrences approximatives de P dans T (moins de k erreurs).

		G	T	C	A	G	G	...
	0	0	0	0	0	0	0	...
C	1	1	1	0	1	1	1	
A	2	2	2	1	0	1	2	
T	3	3	2	2	1	1	2	

Diagram illustrating the dynamic programming table for approximate motif search. The table shows the edit distance (number of errors) between the motif P (rows: C, A, T) and the sequence T (columns: G, T, C, A, G, G, ...). The first row and column are initialized to 0. The values in the table represent the minimum number of errors required to align the motif with the sequence. Red arrows indicate the path from the bottom-right cell (T, G) back to the top-left cell (0, 0), showing the alignment: T aligned with G (1 error), A aligned with A (0 errors), and C aligned with T (1 error). The values 1 and 1 in the bottom row are highlighted with red boxes.

- initialiser la première ligne à 0, et la première colonne comme d'habitude
- même relations de récurrence que pour l'alignement global de séquences
- rechercher à la ligne m toutes les cases contenant des valeurs plus petites ou égales à k
- pour trouver un alignement, suivre les pointeurs jusqu'à la première ligne

Alignement local - Algorithme de Smith-Waterman

Relations de récurrence: $V(0, j) = 0$
 $V(i, 0) = 0$

$$V(i, j) = \max \begin{cases} 0 \\ V(i-1, j) + p(s_i, -) \\ V(i, j-1) + p(-, t_j) \\ V(i-1, j-1) + p(s_i, t_j) \end{cases}$$

- Le 0 dans la récurrence permet d'ignorer un nombre quelconque de caractères en début de séquence
- Pour trouver un alignement local de score maximal:
 - On remplit la table
 - On recherche une case c contenant la valeur maximale de la table
 - De cette case c, on suit les pointeurs jusqu'à une case contenant la valeur 0

Alignement local - Algorithme de Smith-Waterman

$$\begin{aligned}
 V(0, j) &= 0 \\
 V(i, 0) &= 0 \\
 V(i, j) &= \max \begin{cases} 0 \\ V(i-1, j) + p(s_i, -) & -2 \\ V(i, j-1) + p(-, t_j) & -2 \\ V(i-1, j-1) + p(s_i, t_j) & -1 \text{ si } \neq \\ & +2 \text{ si } = \end{cases}
 \end{aligned}$$

Exemple: Calculez la valeur d'un alignement local entre les mots ACGTTA et AACTA

	-	A	C	G	T	T	A
-	0	0	0	0	0	0	0
A	0	2	0	0	0	0	2
A	0	2	1	0	0	0	2
C	0	0	4	2	0	0	0
T	0	0	2	3	4	2	0
A	0	2	0	1	2	3	4

Alignement local - Algorithme de Smith-Waterman

Exemple: Calculez la valeur d'un alignement local entre les mots ACGTTA et AACTA

	-	A	C	G	T	T	A
-	0	0	0	0	0	0	0
A	0	2	0	0	0	0	2
A	0	2	1	0	0	0	2
C	0	0	4	2	0	0	0
T	0	0	2	3	4	2	0
A	0	2	0	1	2	3	4

A C

A C

A C G T

A C - T

T A

T A