



Table des matières

1	Questions de cours	1
1.1	Alignement multiple exact	1
1.2	ClustalW	3
2	Exercice d'alignement	4
EXERCICE 1	Alignement de deux alignements	4
EXERCICE 2	Un autre alignement de deux alignements	7

1 Questions de cours

1.1 Alignement multiple exact

1. Qu'est-ce qu'un alignement multiple?

Correction:

La superposition de plus de 2 séquences

2. Quels sont les 4 types de résultats pouvant être obtenus sur un site dans un alignement multiple de séquences, et à quoi correspondent-ils?

Correction:

- (a) **Les trous** : l'alignement présente une ou plusieurs insertions/délétions (indel) sur un site donné
- (b) **Les correspondances** : tous les résidus présents sur un site sont exactement les mêmes
- (c) **Les différences** : les résidus du site concerné ne sont pas les mêmes et à des fréquences d'apparition différentes (non-équiprobables)
- (d) **Les ambiguïtés** : différents résidus sont présents sur un site donné mais à une fréquence d'apparition équiprobable

3. Rappeler les 2 règles fondamentales à respecter pour qu'un alignement multiple soit valide?

Correction:

- (a) On ne peut pas avoir de superposition composée uniquement de trous (c'est-à-dire

qu'un site ne peut pas contenir que des trous)

(b) Une fois alignées, toutes les séquences doivent avoir la même taille

4. Rappeler et expliquer de quoi se compose un schéma d'évaluation ?

Correction:

Un schéma d'évaluation est composé de :

- (a) la matrice de substitution Σ , qui donne le coût d'une correspondance ou d'une substitution
- (b) la fonction d'évaluation des trous γ , qui donne le coût d'un trou en fonction de sa taille. Cette fonction peut être constante, linéaire ou affine

5. Donner et expliquer les formules permettant d'effectuer un alignement global de 3 séquences, quelle que soit la fonction γ ?

Correction:

Les algorithmes d'alignement multiple sont de simples extensions des algorithmes d'alignement par paire Needleman-Wunsch et Smith-Waterman.

Pour 3 séquences, la matrice de similarité S et la matrice de backtrack B sont désormais des tableaux à 3 dimensions, et la formule de remplissage de l'intérieur de ces cubes est la suivante :

$$s_{i,j,k} = \max \left\{ \begin{array}{l} s_{i-1,j-1,k-1} + \sigma_{x_i,y_j} + \sigma_{y_j,z_k} + \sigma_{x_i,z_k} \longrightarrow \text{Corresp. / sub.} \\ \max(s_{i-r,j,k} + \gamma(x_{(i-r):i}) + \sigma_{y_j,z_k}), r = 1 \dots i \\ \max(s_{i,j-r,k} + \gamma(y_{(j-r):j}) + \sigma_{x_i,z_k}), r = 1 \dots j \\ \max(s_{i,j,k-r} + \gamma(z_{(k-r):k}) + \sigma_{x_i,y_j}), r = 1 \dots k \end{array} \right\} \text{Indel + Corresp. / sub.}$$

$$\left\{ \begin{array}{l} \max(s_{i-r,j-l,k} + \gamma(x_{(i-r):i}) + \gamma(y_{(j-l):j})), r = 1 \dots i, l = 1 \dots j \\ \max(s_{i,j-r,k-l} + \gamma(y_{(j-r):j}) + \gamma(z_{(k-l):k})), r = 1 \dots j, l = 1 \dots k \\ \max(s_{i-r,j,k-l} + \gamma(x_{(i-r):i}) + \gamma(z_{(k-l):k})), r = 1 \dots i, l = 1 \dots k \end{array} \right\} \text{Indel} \times 2$$

Pour les bordures, les formules sont également similaires à celles pour un alignement par paire :

$$\begin{aligned} s_{i,0,0} &= \gamma(x_{1:i}) \\ s_{0,j,0} &= \gamma(y_{1:j}) \\ s_{0,0,k} &= \gamma(z_{1:k}) \end{aligned}$$

6. Quelle est la formule permettant d'exprimer en fonction de n le nombre de cas à considérer (nombre de lignes) pour le remplissage de l'intérieur de la matrice S ?

Correction:

$$\text{Formule : } 1 + \sum_{k=1}^{n-1} \binom{n}{k}$$

7. De quel ordre sont les complexités temporelles et spatiales des algorithmes d'alignement multiple exacts ? Que peut-on en dire sur l'efficacité de tels algorithmes ?

Correction:

Les complexités spatiales et temporelles sont en $O(x^n)$, c'est-à-dire qu'elles sont de l'ordre de x^n , ce qui correspond à une complexité exponentielle.

Elles indiquent donc que les algorithmes d'alignement multiple exact ne sont pas effi-

caces, c'est pourquoi d'autres méthodes ont été développées.

1.2 ClustalW

1. Qu'est-ce que la méthode de résolution par force brute ?

Correction:

Méthode consistant à tester toutes les valeurs d'entrée possibles de la fonction f que l'on souhaite optimiser.

Elle peut s'appliquer à tous les problèmes mais nécessite généralement beaucoup de calculs.

2. La méthode d'alignement multiple exacte est-elle une méthode par force brute ? Si non, de quel type de méthode s'agit-il ?

Correction:

Les algorithmes utilisés pour l'alignement multiple exact sont des extensions de Needleman-Wunsh et Smith-Waterman, qui sont des algorithmes de résolution par force brute. La méthode d'alignement multiple exacte est donc une méthode par force brute.

3. Qu'est-ce qu'une méthode de résolution par heuristique ?

Correction:

Méthode consistant à ajouter une ou plusieurs hypothèses à un problème afin de le simplifier. Ainsi il est possible de le résoudre plus facilement.

La solution ainsi trouvée est valide mais il n'est pas garanti qu'elle soit optimale, c'est-à-dire que ce soit la meilleure.

4. Qu'est-ce que l'heuristique d'"alignement progressif" ?

Correction:

Hypothèse selon laquelle il est possible d'obtenir un bon alignement multiple à partir d'alignements par paires successifs. L'ordre de ces alignements est donné par l'"arbre guide".

5. Comment définiriez-vous un arbre binaire ?

Correction:

Structure de donnée récursive représentée sous forme de nœuds reliés entre eux.

Un arbre peut être vide. Dans le cas contraire, il contient un nœud avec une valeur, et ce nœud possède 0, 1 ou 2 fils (gauche et droit), qui sont eux aussi des arbres.

La racine de l'arbre représente le nœud n'ayant pas de parent, et les feuilles sont les nœuds dont les 2 fils sont des arbres vides.

6. Qu'est-ce que le profil d'un alignement ?

Correction:

Matrice indiquant la fréquence d'apparition de chaque caractère ou de trou pour chaque colonne de l'alignement.



7. Rappeler et expliquer la formule permettant de remplir une matrice de substitution entre 2 alignements.

Correction:

$$\text{Formule : } \sigma(P_A, P_B)_{i,j} = \sum_{k=1}^{|R|} \sum_{l=1}^{|R|} p_{A_{k,i}} \times p_{B_{l,j}} \times \sigma_{k,l}$$

Chaque case de la matrice contient le score d'une substitution de la i ème colonne de A vers la j ème colonne de B , correspondant à la somme des produits de chaque paire de résidus pondérés par leur score de substitution.

8. Écrire et expliquer les complexités temporelles et spatiales de l'algorithme ClustalW, pour n séquences. Faire ce travail pour le cas général, et celui où la fonction γ est linéaire

Correction:

Se référer aux explications données lors du cours. A noter que les complexités temporelles sont linéaires en n tandis que les complexités spatiales sont constantes ce qui est très performant !

2 Exercice d'alignement

Exercice 1: Alignement de deux alignements

Soient les 2 alignements suivants que nous souhaitons aligner :

Alignement A :

TAACG-
TTACCA

Alignement B :

TGCTCT
TCAT-G

Nous utiliserons pour cela l'algorithme de Needleman-Wunsch global avec le schéma d'évaluation suivant :

	A	T	G	C	-
A	1	0	0.5	0	0
T	0	1	0	0.5	0
G	0.5	0	1	0	0
C	0	0.5	0	1	0
-	0	0	0	0	0

$\Sigma =$, $\gamma(g) = -1 \times |g|$

1. Commencez par calculer les profils P_A et P_B :

Correction:

P_A	T	A	A	C	G	-	P_B	T	G	C	T	C	T
	T	T	A	C	C	A		T	C	A	T	-	G
A		1/2	1			1/2	A			1/2			
T	1	1/2					T	1			1		1/2
C				1	1/2		C		1/2	1/2		1/2	
G					1/2		G		1/2				1/2
-						1/2	-					1/2	

2. Complétez ensuite les case vides de la matrice de substitution des profils $\Sigma(P_A, P_B)$:

Correction:

		T	G	C	T	C	T
		T	C	A	T	-	G
T	T	1	1/4	1/4	1	1/4	1/2
A	T	1/2	1/4	3/8	1/2	1/8	3/8
A	A	0	1/4	1/2	0	0	1/4
C	C	1/2	1/2	1/2	1/2	1/2	1/4
G	C	1/4	1/2	3/8	1/4	1/4	3/8
-	A	0	1/8	1/4	0	0	1/8

3. Puis complétez les matrices de similarité et de backtrack ci-dessous :

Correction:

S		∅	T	G	C	T	C	T
		∅	T	C	A	T	-	G
∅	∅	0	-1	-2	-3	-4	-5	-6
T	T	-1	1	0	-1	-2	-3	-4
A	T	-2	0	5/4	3/8	-1/2	-3/2	-5/2
A	A	-3	-1	1/4	7/4	3/4	-1/4	-5/4
C	C	-4	-2	-1/2	3/4	9/4	5/4	1/4
G	C	-5	-3	-3/2	-1/8	5/4	5/2	13/8
-	A	-6	-4	-5/2	-9/8	1/4	3/2	21/8

B		∅	T	G	C	T	C	T
		∅	T	C	A	T	-	G
∅	∅	∅	←	←	←	←	←	←
T	T	↑	↖	←	←	↖	←	←
A	T	↑	↑	↖	↖	↖	←	←
A	A	↑	↑	↖	↖	←	←	↖
C	C	↑	↑	↖	↖	↖	↖	←
G	C	↑	↑	↖	↖	↑	↖	↖
-	A	↑	↑	↑	↑	↑	↑	↖

4. Enfin écrivez l'alignement trouvé dans le tableau suivant :

Correction:

A :	TAACG-
	TTACCA
B :	TGCTCT
	TCAT-G

Exercice 2: Un autre alignement de deux alignements

Aligner les 2 alignements ci-dessous en suivant l'algorithme de Needleman-Wunsh semi-global :

A :
CTACA
CTATG

B :
CTATCA
CTAGTA
CTAGCA

Vous utiliserez pour cela la matrice de substitution Σ et la fonction d'évaluation des trous γ suivantes :

$$\Sigma = \begin{array}{c|ccccc} & A & T & G & C & - \\ \hline A & 1 & 0 & 0.5 & 0 & 0 \\ T & 0 & 1 & 0 & 0.5 & 0 \\ G & 0.5 & 0 & 1 & 0 & 0 \\ C & 0 & 0.5 & 0 & 1 & 0 \\ - & 0 & 0 & 0 & 0 & 0 \end{array}, \gamma(g) = -0,5 \times |g|$$

Correction:

Comme indiqué dans le cours, il y a 3 étapes à réaliser pour obtenir l'alignement :

1. Calculer les profils P_A et P_B :

P_A	C	T	A	C	A
	C	T	A	T	G
A			1		1/2
T		1		1/2	
C	1			1/2	
G					1/2
-					

P_B	C	T	A	T	C	A
	C	T	A	G	T	A
	C	T	A	G	C	A
A			1			1
T		1		1/3	1/3	
C	1				2/3	
G				2/3		
-						

2. Calculer la matrice de substitution des profils $\Sigma(P_A, P_B)$:

On doit utiliser la formule du cours pour remplir chaque case de la matrice, en ne s'intéressant qu'aux paires de résidus dont les valeurs dans les profils P_A et P_B ne sont pas nulles (ce qui est fait dans la correction ci-dessous).

La matrice finale est la suivante (contrairement aux profils, il faut bien écrire les 0) :

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

Vous trouverez ci-dessous les formules détaillées pour chaque case :

$$\begin{aligned}\sigma(P_A, P_B)_{0,0} &= p_{A_{C,0}} \times p_{B_{C,0}} \times \sigma_{C,C} \\ &= 1 \times 1 \times 1 \\ &= 1\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{0,1} &= p_{A_{C,0}} \times p_{B_{T,1}} \times \sigma_{C,T} \\ &= 1 \times 1 \times 1/2 \\ &= 1/2\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{0,2} &= p_{A_{C,0}} \times p_{B_{A,2}} \times \sigma_{C,A} \\ &= 1 \times 1 \times 0 \\ &= 0\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{0,3} &= p_{A_{C,0}} \times p_{B_{T,3}} \times \sigma_{C,T} + p_{A_{C,0}} \times p_{B_{G,3}} \times \sigma_{C,G} \\ &= 1 \times 1/3 \times 1/2 + 1 \times 2/3 \times 0 \\ &= 1/6\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{0,4} &= p_{A_{C,0}} \times p_{B_{T,4}} \times \sigma_{C,T} + p_{A_{C,0}} \times p_{B_{C,4}} \times \sigma_{C,C} \\ &= 1 \times 1/3 \times 1/2 + 1 \times 2/3 \times 1 \\ &= 5/6\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{0,5} &= p_{A_{C,0}} \times p_{B_{A,5}} \times \sigma_{C,A} \\ &= 1 \times 1 \times 0 \\ &= 0\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{1,0} &= p_{A_{T,1}} \times p_{B_{C,0}} \times \sigma_{T,C} \\ &= 1 \times 1 \times 1/2 \\ &= 1/2\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{1,1} &= p_{A_{T,1}} \times p_{B_{T,1}} \times \sigma_{T,T} \\ &= 1 \times 1 \times 1 \\ &= 1\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{1,2} &= p_{A_{T,1}} \times p_{B_{A,2}} \times \sigma_{T,A} \\ &= 1 \times 1 \times 0 \\ &= 0\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{1,3} &= p_{A_{T,1}} \times p_{B_{T,3}} \times \sigma_{T,T} + p_{A_{T,1}} \times p_{B_{G,3}} \times \sigma_{T,G} \\ &= 1 \times 1/3 \times 1 + 1 \times 2/3 \times 0 \\ &= 1/3\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{1,4} &= p_{A_{T,1}} \times p_{B_{T,4}} \times \sigma_{T,T} + p_{A_{T,1}} \times p_{B_{C,4}} \times \sigma_{T,C} \\ &= 1 \times 1/3 \times 1 + 1 \times 2/3 \times 1/2 \\ &= 2/3\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{1,5} &= p_{A_{T,1}} \times p_{B_{A,5}} \times \sigma_{T,A} \\ &= 1 \times 1 \times 0 \\ &= 0\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{2,0} &= p_{A_{A,2}} \times p_{B_{C,0}} \times \sigma_{A,C} \\ &= 1 \times 1 \times 0 \\ &= 0\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{2,1} &= p_{A_{A,2}} \times p_{B_{T,1}} \times \sigma_{A,T} \\ &= 1 \times 1 \times 0 \\ &= 0\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{2,2} &= p_{A_{A,2}} \times p_{B_{A,2}} \times \sigma_{A,A} \\ &= 1 \times 1 \times 1 \\ &= 1\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{2,3} &= p_{A_{A,2}} \times p_{B_{T,3}} \times \sigma_{A,T} + p_{A_{A,2}} \times p_{B_{G,3}} \times \sigma_{A,G} \\ &= 1 \times 1/3 \times 0 + 1 \times 2/3 \times 1/2 \\ &= 1/3\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{2,4} &= p_{A_{A,2}} \times p_{B_{T,4}} \times \sigma_{A,T} + p_{A_{A,2}} \times p_{B_{C,4}} \times \sigma_{A,C} \\ &= 1 \times 1/3 \times 0 + 1 \times 2/3 \times 0 \\ &= 0\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{2,5} &= p_{A_{A,2}} \times p_{B_{A,5}} \times \sigma_{A,A} \\ &= 1 \times 1 \times 1 \\ &= 1\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{3,0} &= p_{A_{T,3}} \times p_{B_{C,0}} \times \sigma_{T,C} + p_{A_{C,3}} \times p_{B_{C,0}} \times \sigma_{C,C} \\ &= 1/2 \times 1 \times 1/2 + 1/2 \times 1 \times 1 \\ &= 3/4\end{aligned}$$

$$\begin{aligned}\sigma(P_A, P_B)_{3,1} &= p_{A_{T,3}} \times p_{B_{T,1}} \times \sigma_{T,T} + p_{A_{C,3}} \times p_{B_{T,1}} \times \sigma_{C,T} \\ &= 1/2 \times 1 \times 1 + 1/2 \times 1 \times 1/2 \\ &= 3/4\end{aligned}$$

$$\begin{aligned}
\sigma(P_A, P_B)_{3,2} &= p_{A_{T,3}} \times p_{B_{A,2}} \times \sigma_{T,A} + p_{A_{C,3}} \times p_{B_{A,2}} \times \sigma_{C,A} \\
&= 1/2 \times 1 \times 0 + 1/2 \times 1 \times 0 \\
&= 0 \\
\sigma(P_A, P_B)_{3,3} &= p_{A_{T,3}} \times p_{B_{T,3}} \times \sigma_{T,T} + p_{A_{T,3}} \times p_{B_{G,3}} \times \sigma_{T,G} + p_{A_{C,3}} \times p_{B_{T,3}} \times \sigma_{C,T} + \\
&\quad p_{A_{C,3}} \times p_{B_{G,3}} \times \sigma_{C,G} \\
&= 1/2 \times 1/3 \times 1 + 1/2 \times 2/3 \times 0 + 1/2 \times 1/3 \times 1/2 + 1/2 \times 2/3 \times 0 \\
&= 1/4 \\
\sigma(P_A, P_B)_{3,4} &= p_{A_{T,3}} \times p_{B_{T,4}} \times \sigma_{T,T} + p_{A_{T,3}} \times p_{B_{C,4}} \times \sigma_{T,C} + p_{A_{C,3}} \times p_{B_{T,4}} \times \sigma_{C,T} + \\
&\quad p_{A_{C,3}} \times p_{B_{C,4}} \times \sigma_{C,C} \\
&= 1/2 \times 1/3 \times 1 + 1/2 \times 2/3 \times 1/2 + 1/2 \times 1/3 \times 1/2 + 1/2 \times 2/3 \times 1 \\
&= 3/4 \\
\sigma(P_A, P_B)_{3,5} &= p_{A_{T,3}} \times p_{B_{A,5}} \times \sigma_{T,A} + p_{A_{C,3}} \times p_{B_{A,5}} \times \sigma_{C,A} \\
&= 1/2 \times 1 \times 0 + 1/2 \times 1 \times 0 \\
&= 0 \\
\sigma(P_A, P_B)_{4,0} &= p_{A_{A,4}} \times p_{B_{C,0}} \times \sigma_{A,C} + p_{A_{G,4}} \times p_{B_{C,0}} \times \sigma_{G,C} \\
&= 1/2 \times 1 \times 0 + 1/2 \times 1 \times 0 \\
&= 0 \\
\sigma(P_A, P_B)_{4,1} &= p_{A_{A,4}} \times p_{B_{T,1}} \times \sigma_{A,T} + p_{A_{G,4}} \times p_{B_{T,1}} \times \sigma_{G,T} \\
&= 1/2 \times 1 \times 0 + 1/2 \times 1 \times 0 \\
&= 0 \\
\sigma(P_A, P_B)_{4,2} &= p_{A_{A,4}} \times p_{B_{A,2}} \times \sigma_{A,A} + p_{A_{G,4}} \times p_{B_{A,2}} \times \sigma_{G,A} \\
&= 1/2 \times 1 \times 1 + 1/2 \times 1 \times 1/2 \\
&= 3/4 \\
\sigma(P_A, P_B)_{4,3} &= p_{A_{A,4}} \times p_{B_{T,3}} \times \sigma_{A,T} + p_{A_{A,4}} \times p_{B_{G,3}} \times \sigma_{A,G} + p_{A_{G,4}} \times p_{B_{T,3}} \times \sigma_{G,T} + \\
&\quad p_{A_{G,4}} \times p_{B_{G,3}} \times \sigma_{G,G} \\
&= 1/2 \times 1/3 \times 0 + 1/2 \times 2/3 \times 1/2 + 1/2 \times 1/3 \times 0 + 1/2 \times 2/3 \times 1 \\
&= 1/2 \\
\sigma(P_A, P_B)_{4,4} &= p_{A_{A,4}} \times p_{B_{T,4}} \times \sigma_{A,T} + p_{A_{A,4}} \times p_{B_{C,4}} \times \sigma_{A,C} + p_{A_{G,4}} \times p_{B_{T,4}} \times \sigma_{G,T} + \\
&\quad p_{A_{G,4}} \times p_{B_{C,4}} \times \sigma_{G,C} \\
&= 1/2 \times 1/3 \times 0 + 1/2 \times 2/3 \times 0 + 1/2 \times 1/3 \times 0 + 1/2 \times 2/3 \times 0 \\
&= 0 \\
\sigma(P_A, P_B)_{4,5} &= p_{A_{A,4}} \times p_{B_{A,5}} \times \sigma_{A,A} + p_{A_{G,4}} \times p_{B_{A,5}} \times \sigma_{G,A} \\
&= 1/2 \times 1 \times 1 + 1/2 \times 1 \times 1/2 \\
&= 3/4
\end{aligned}$$

3. Réaliser l'alignement de A et B :

On remplit les matrices S et B comme pour un alignement par paires :

S		∅	C	T	A	T	C	A
	∅	∅	C	T	A	G	T	A
	∅	∅	C	T	A	G	C	A
∅	∅	0	0	0	0	0	0	0
C	C	0	1	0.5	0	0.16	0.83	0.33
T	T	0	0.5	2	1.5	1	0.82	0.83
A	A	0	0	1.5	3	2.5	2	1.82
C	T	0	0.75	1	2.5	3.25	3.25	2.75
A	G	0	0.25	0.75	2	3	3.25	4

B		∅	C	T	A	T	C	A
	∅	∅	C	T	A	G	T	A
	∅	∅	C	T	A	G	C	A
∅	∅	∅	∅	∅	∅	∅	∅	∅
C	C	∅	↖	↖	↖	↖	↖	←
T	T	∅	↖	↖	←	←	↖	↖
A	A	∅	↖	↑	↖	←	←	↖
C	T	∅	↖	↑	↑	↖	↖	←
A	G	∅	↑	↖	↑	↖	↖	↖

On identifie le maximum parmi les bordures du bas et de droite de la matrice S (car alignement semi-global), puis on remonte le chemin de la matrice B pour trouver l'alignement :

A : CTA-CA
 CTA-TG

 CTATCA
 B : CTAGTA
 CTAGCA

NB : vous pouvez écrire A au-dessus ou en-dessous de B, ça n'a pas d'importance ici (alignement d'alignements)