# TD 2 : Consensus, Motif, Profil

Anne-Muriel Arigon Chifolleau - Université Montpellier, LIRMM  ${\it chifolleau@lirmm.fr}$ 

#### 1 Consensus

Code	Description
A	Adenine
С	Cytosine
G	Guanine
$\mathbf{T}$	Thymine
U	Uracil
R	Purine (A or G)
Y	Pyrimidine (C, T, or U)
M	C or A
K	T, U, or G
W	T, U, or A
S	C or G
В	C, T, U, or G (not A)
D	A, T, U, or G (not C)
Н	A, T, U, or C (not G)
V	A, C, or G (not T, not U)
N	Any base (A, C, G, T, or U)

Voici ci-contre, le code IUPAC (International Union of Pure and Applied Chemistry) pour les nucléotides.

Une séquence **consensus** est une séquence pour laquelle à chaque position de l'alignement multiple, on retient la lettre majoritaire.

Ci-dessous, un alignement multiple du site de fixation du facteur de transcription c-Ets-1 chez les 15 murins (source TRANSFAC M00032).

Trouver le consensus pour ce site de fixation ainsi que son code IUPAC.

	G	$^{\circ}$	С	G	G	A	A	G	T	G
	A	С	С	G	G	A	A	G	С	A
	G	С	С	G	G	A	Т	G	T	A
	A	С	С	G	G	A	A	G	С	T
	A	С	С	G	G	A	Т	A	T	A
	С	С	С	G	G	A	A	G	T	G
	A	С	Α	G	G	A	A	G	T	С
	G	С	С	G	G	A	Т	G	С	A
	Т	С	С	G	G	A	A	G	T	A
	A	С	Α	G	G	A	A	G	С	G
	A	С	Α	G	G	A	Т	A	T	G
	Т	С	С	G	G	A	A	A	С	С
	A	С	Α	G	G	A	$\Gamma$	A	T	С
	С	Α	Α	G	G	A	С	G	Α	С
	Т	С	Τ	G	G	A	С	С	С	T
Séquence consensus :										
-										
Code IUPAC :										

## 2 Représentation d'un motif par une matrice

## 2.1 Comptage

Dans la matrice que nous voulons remplir, les lignes représentent les colonnes de l'alignement multiple et les colonnes de la matrice, les acides nucléiques (4 colonnes).

Avec l'exemple de l'exercice précédent, remplissez la matrice ci-dessous. La première ligne est donnée à titre d'exemple.

A	C 2	G 3	T 3
7	2	3	3

## 2.2 Matrice de fréquence et matrice poids position

Une fois le comptage effectué, on détermine la fréquence de chaque nucléotides dans chaque colonne du MSA (valeur entre 0 et 1). Pour passer à une matrice poids-position il faut encore une opération supplémentaire. L'idée est que les bases qui apparaissent plus que la moyenne aient un score positif et celles qui apparaissent moins que la moyenne, un score négatif. On applique donc la formule suivante où f(x) représente la fréquence de x.

Poids de la base x dans une colonne de l'alignement :  $\log \frac{f(x)}{0.25}$ 

Que représente 0.25?

Les matrices de fréquences et poids-position correspondant à notre exemple sont données ci-dessous.

	A	С	G	Т
0	.47	0.13	0.2	0.2
0	.07	0.93	0	0
0	.33	0.6	0	0.07
	0	0	1	0
	0	0	1	0
	1	0	0	0
0	.53	0.13	0	0.33
0	.27	0.07	0.67	0
0	.07	0.4	0	0.53
0	.33	0.27	0.27	0.13

Matrice de fréquences

A	С	G	T
0.91	-0.94	-0.32	-0.32
-1.8	1.9	-2.3	-2.3
0.4	1.26	-2.3	-1.8
-2.3	-2.3	2	-2.3
-2.3	-2.3	2	-2.3
2	-2.3	-2.3	-2.3
1.1	-0.94	-2.3	0.4
0.11	-1.8	1.42	-2.3
-1.8	0.68	-2.3	1.1
0.4	0.11	0.11	-0.94

Matrice poids-position

#### 2.3 Calcul du score d'un mot

Une fois la matrice poids-position calculée. On s'intéresse au score d'un mot dans ce modèle. Pour calculer le score d'un mot, il suffit d'additionner les poids correspondant aux lettres qui le compose, position par position. Considérons la matrice poids-position de la feuille précédente.

- 1. Comment sait-on qu'un score sera bon ou mauvais?
- 2. Calculer le meilleur score possible avec cette matrice; quel motif atteint ce score?
- 3. Le moins bon score; de même, quel motif atteint ce score?
- 4. Calculer le score du mot TACGGATACG

#### 3 Motif PROSITE

Nous allons maintenant créer un motif PROSITE correspondant à l'alignement multiple suivant d'hormone pancréatique. (Source : PROSITE, PS00265)

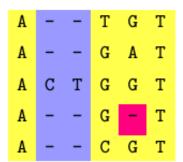
```
NEUY CARAU/29-64 AEE..LAKYYSALRHYINLITRQRY
             PEE. . L NRYYASLRHYLNL VTRQRY
PYY HUMAN/29-64
PMY PETMA/1-36
             P E E . . L S K Y M L A V R N Y I N L I T R Q R Y
             PED. . WASYQAAVRHYVNLITRQRY
PPY LOPAM/1-36
             PEQ. . . MAQYAAELRRYINMLTRPRY
PAHO BOVIN/30-65
             V E D . . L I R F Y N D L Q Q Y L N V V T R H R Y
PAHO CHICK/26-61
PAHO ANSAN/1-36
             VED. . LRFYYDNLQQYRLNVFRHRY
             P N E . . L R Q Y L K E L N E Y Y A I M G R T R F
NPF HELAS/4-39
NPF MONEX/1-39
             DNKAALRDYLRQINEYFAI I GRPRF
```

La syntaxe pour la description du motif est la suivante :

- Tous les éléments de l'expression sont séparés par des tirets -
- Le joker est la lettre X pour les protéines (N pour l'ADN)
- On peut préciser le nombre d'occurrences avec des parenthèses : X(5) ou D(2,4)
- Le choix entre plusieurs acides aminés possibles se note avec des crochets [AP]
- 1. À partir de l'alignement multiple ci-dessous, construisez un motif pour le site actif présumé.
- 2. Combien de "chaînes de caractères" correspondent au motif que vous venez de trouver?
- 3. Combien de "chaînes de caractères" correspondent au motif 'D-E-X(3)-[LIVM]-X(1,3)-[FY] '?

# 4 Modèle de Markov caché (HMM)

Vous devez construire les profils HMM des 2 alignements multiples ci-dessous .



G	Α	C	С	Α
G	Α	C	С	Α
G	T	G	С	Α
G	Α	C	С	T
G	Т	G	С	G