



SORBONNE UNIVERSITÉS

3I005 - STATISTIQUE ET INFORMATIQUE

PROJET 2 - ANALYSE STATISTIQUE D'UNE FAMILLE DE PROTEINES

PAR

MARC HU & SABRINA CHAOUCHE

Année universitaire : 2017-2018

Table des matières

I	Réponses aux questions	2
II	Guide d'utilisation	4

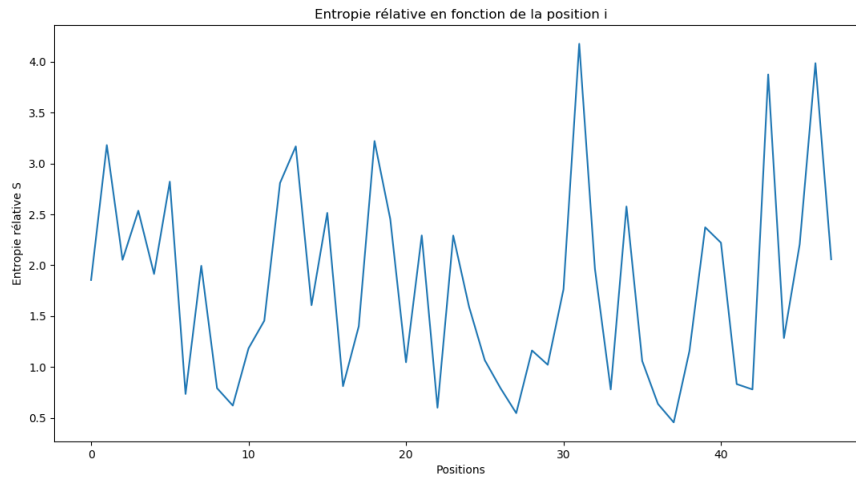
I Réponses aux questions

I.1 Données

Les données des fichiers Dtrain, test-seq, test-seq2 et distances sont lues et conservées dans des matrices pour leurs utilisations.

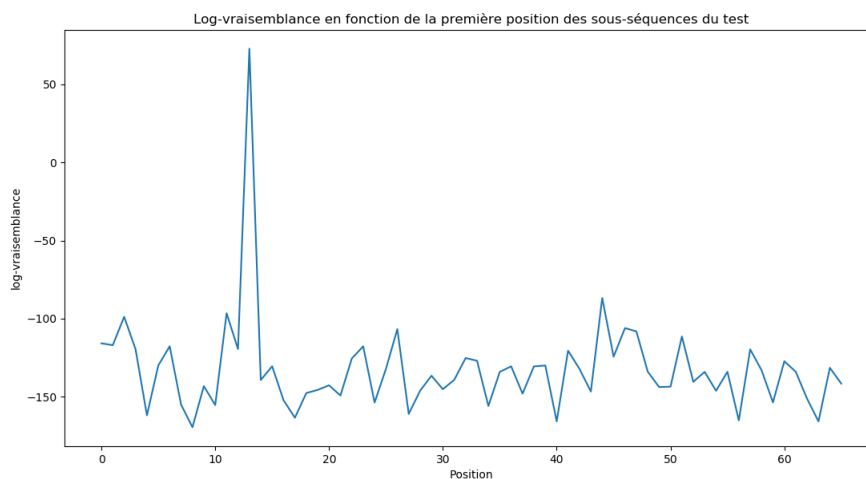
I.2 Modélisation par PSWM

L'entropie relative S :



La valeur de l'entropie pour une position donnée "i" varie entre 0 et $\log_2(21)$: plus la valeur est maximale, plus la position est conservée. D'après le tracé de l'entropie, on remarque que la position la plus conservée est la position 31 et l'acide la plus conservée à cette position est l'acide W.

La log-vraisemblance :

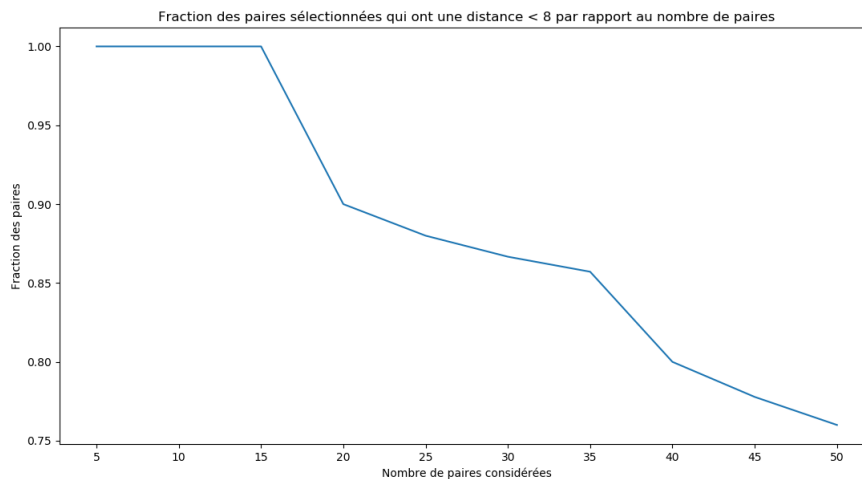


La log-vraisemblance indique l'appartenance d'une séquence à la famille donnée par Dtrain si sa valeur est positive. D'après le tracé de la log-vraisemblance des sous-séquences du fichier de test, on a une seule valeur positive donc on en déduit qu'il y a une seule séquence qui appartient à la famille donnée par Dtrain.

La sous-séquence du fichier de test appartenant à la famille donnée par Dtrain est :
KALYDFAGQSAGELSLGKDEIILVTQKENNGWWLASRLDKSASGWAPS.

I.3 Co-évolution des résidues en contact :

La fraction des 50 paires sélectionnées



Le tracé de la fraction des 50 paires sélectionnées montre que plus les paires sont corrélées plus la probabilité qu'elles soient en contact augmente : plus l'information mutuelle M_{ij} augmente plus la probabilité que les positions correspondantes soient statistiquement dépendantes. Pour les 10 plus grandes valeurs M_{ij} , l'information mutuelle est très élevée donc leurs distances respectives sont < 8 et dans ce cas, on a une *fraction* = 1 ce qui signifie que les paires sont en contact. Dans le cas où on considère les 50, la fraction diminue car l'information mutuelle diminue.

II Guide d'utilisation

Pour faciliter le test des fonctions, on a utilisé un menu depuis la ligne de commande. Pour lancer les menu, il suffit de taper la commande **"python3 projet2.py"** et le menu suivant s'affiche : pour utiliser les fonctions de la partie "Modélisation par PSWM" il suffit d'entrer **1** et pour les fonctions de la partie "Co-évolution des résidues en contact", entrer **2**. Pour quitter le programme, entrez **0**.

```
Projet 3I005,  
Statistique en Bioinformatique,  
Analyse statique d'une famille de protéines,  
Par Chaouche Sabrina & Hu Marc  
Choisissez la partie du projet que vous souhaitez:  
1. Première partie  
2. Deuxième partie  
0. Quitter  
>> |
```

La première partie :

Pour utiliser une des fonctions de la partie "Modélisation par PSWM", il suffit de choisir un argument entre 1 et 4 correspondant au numéro de la fonction -tel est mentionné dans l'énoncé- et de suivre les instructions pour effectuer le calcul correspondant.

```
Première partie  
1. Fonction 1 (wi(a))  
2. Fonction 2 (nij(a, b) et wij(a, b))  
3. Fonction 3 (Information mutuelle Mij)  
4. Fonction 4 (Courbe fraction selon le nombre de paires)  
9. Retour  
0. Quitter  
>> |
```

Exemple d'utilisation de la fonction 2 de la partie "Modélisation par PSWM"

```
Première partie  
1. Fonction 1 (ni(a), wi(a))  
2. Fonction 2 (Entropie)  
3. Fonction 3 (Paramètre f(o)(a))  
4. Fonction 4 (Log-vraisemblance)  
9. Retour  
0. Quitter  
>> 2  
En attente de calcul pour la matrice d'entrainement et les wia...  
Quelle position?  
>> 0  
Les trois positions les plus conservées : [ 31. 46. 43.]  
Acides les plus conservés aux positions les plus conservés : ['W', 'P', 'G']  
Si à la position 0 est de 1.85476633369  
  
Appuyer sur une touche pour revenir au menu  
|
```

La seconde partie :

Pour utiliser une des fonctions de la partie "Co-évolution des résidues en contact", il suffit de choisir un argument entre 1 et 4 correspondant au numéro de la fonction -tel est mentionné dans l'énoncé- et de suivre les instructions pour effectuer le calcul correspondant.

```
Deuxième partie
1. Fonction 1 (wi(a))
2. Fonction 2 (nij(a, b) et wij(a, b))
3. Fonction 3 (Information mutuelle Mij)
4. Fonction 4 (Courbe fraction selon le nombre de paires)

9. Retour
0. Quitter
>> |
```

Exemple d'utilisation de la fonction 4 de la partie "Co-évolution des résidues en contact"

```
1. Fonction 1 (wi(a))
2. Fonction 2 (nij(a, b) et wij(a, b))
3. Fonction 3 (Information mutuelle Mij)
4. Fonction 4 (Courbe fraction selon le nombre de paires)

9. Retour
0. Quitter
>> 3
En attente de calcul pour la matrice d'entraînement et wia...

Quelle position i?
>> 0
Quelle position j (strictement supérieur à i)?
>> 1
Calcul de Wij...

Calcul de l'information mutuelle Mij...

L'information mutuelle pour i = 0 et j = 1 est 0.404047499991

Appuyer sur une touche pour revenir au menu
|
```