

Présentation de projet

analyse de données biologiques

Bouaoud Malik

Kenai Selma

Université de Lille/Informatique
UE Bio-informatique

12/02/2021

Introduction

- Analyse de données biologiques:
 - Récupération de séquence biologique depuis le NCBI.
 - Extraction d'information.
 - Traitement d'information.
 - Recherche des sites de facteurs de transcription.



Problème :

Rechercher les TF qui régulent un ensemble de gènes supposés co-régulés, grâce aux occurrence de son TFBS.

Solution triviale:

Pour une première solution il suffit d'observer les occurrences d'un TFBS pour un TF donné dans un ensemble de séquences d'ADN en amont du gène.

Et donc on devra procéder par **fenêtre glissante** en récupérant la fenêtre avec le nombre d'occurrence de notre TFBS, le score dans notre cas d'étude.

Solution proposée :

Au lieu de compter simplement les occurrences d'une fenêtre, nous proposons de calculer une moyenne pour chacune des fenêtres. Cette solution nous a paru plus pertinente puisqu'ici nous privilégions beaucoup plus le score de chaque séquence que son nombre d'occurrence.

Étapes :

Cela se fait en 2 étapes principales :

- Traitement d'ensemble de séquences biologiques.
- Personnalisation du calcul du meilleur score.

Traitement d'ensemble de séquences biologiques

1. Récupérer les séquences en amont du gène depuis le NCBI, grâce à **download promoters**
2. utiliser **scan seq** et **scan all sequences** pour récupérer les fenêtres dépassant un certain seuil.
3. Stocker les résultats satisfaisant notre contrainte dans une structure de donnée imbriquée, c'est à dire : **List[tuple(List,identifiant de la sequence)]**

Personalisation du calcul du meilleur score

Algorithm 1: BestWindowPersonalized

Result: (startPosition,EndPosition,BestScore)
Input: (ListOFSequences,WindowsSize,seqSize)
startCoordinate $\leftarrow 0$
endCoordinate $\leftarrow windowSize$
List $\leftarrow []$
For $i \in List$
Moyenne $\leftarrow 0$
if startCoordinate $\neq seqSize$ **and** endCoordinate $\neq seqSize$ **then**
 Moyenne $\leftarrow scanSequence(i[0], startCoordinate, EndCoordinate)$
 List $\leftarrow \mathbf{Append}(StartCoordinate, EndCoordinate, Moyenne)$
 startCoordinate $\leftarrow EndCoordinate$
 EndCoordinate $\leftarrow EndCoordinate + windowSize$
else
end
//car la moyenne est le troisième paramètre de la liste MaxMoy
 $\leftarrow Max(list[3])$
return (x,y,MaxMoy)



Des questions ?