

Laborator 6 – Bioinformatică

I. Alinierea de secvențe

1. Definirea problemei

Alinierea de secvențe reprezintă găsirea subsecvențelor similare dintre două sau mai multe secvențe. Alinierea se realizează utilizând tehnici de potrivire de șiruri de caractere ("*string matching*"). Potrivirea de șiruri de caractere nu este o problemă simplă și este utilizată foarte des în bioinformatică. Astfel, găsirea algoritmilor eficienți și optimi pentru alinieri este foarte importantă.

Există algoritmi complecși pentru alinieri, bazați pe tehnici de programare dinamică sau alte tehnici dificile de programare. Dar algoritmi utilizați diferă în funcție de tipurile de secvențe care trebuie alinate și de tipurile de alinieri.

Aplicațiile alinierilor în bioinformatică sunt numeroase, precum:

1. determinarea relațiilor filogenetice între specii (clasificarea evolutivă);
2. clasificarea proteinelor și stabilirea anumitor funcții impuse de structură;
3. căutarea anumitor gene în cadrul unei secvențe de ADN;
4. compararea anumitor secvențe de ADN (criminalistică);
5. studierea materialelor genetice virale sau bacteriene în scopul găsirii de medicamente adecvate.

Există mai multe tipuri de alinieri:

1. Aliniere globală

Reprezintă alinierea completă a unei secvențe cu o altă secvență, de exemplu:

```
and.--e-.can.sequence--.very.well.now
|||| | |||| | | ||||| |
and.they.can.lear-n.too.very.soo--n--
```

Există mai multe modalități de potrivire globală a celor două șiruri. Se urmărește găsirea celor mai bune potriviri, cu cât mai puține goluri, "gaps" (–).

2. Aliniere locală

Reprezintă găsirea unei regiuni dintr-o secvență care se potrivește cu o regiune din altă secvență, de exemplu:

```

we.try.to.help.people
||||||| |
you.try.to.learn.bioinformatics

```

Exemplu de alinieri diferite pentru două şiruri:

Globală:

```

GTGTAKIKKAVAV
G--TAK-KKA-AV

```

Locală:

```

GTGTAKIKK-AVAV
--GTAK-KKAAV--

```

3. Găsire de secvenţe scurte într-o secvenţă lungă:

```

      match
      |||||
we.shall.match.the.words

```

4. Aliniere multiplă între mai multe secvenţe

Scorul unei alinieri

Aşa cum se poate observa, există mai multe alinieri posibile pentru aceleaşi secvenţe. Scorul unei alinieri reprezintă o valoare care estimează cât de bună este o aliniere. Pentru a alege cea mai bună aliniere trebuie să se țină cont de semnificaţia biologică a secvenţelor şi de proprietăţile structurale ale elementelor componente. Astfel, în urma datelor experimentale şi statistice, s-au definit matrici de scoruri pentru alinieri. În calcularea scorului final al unei alinieri, potrivirile de caractere pot avea ponderi diferite sau egale, în funcţie de problemă.

Exemplu:

Un exemplu de matrice de scoruri pentru alinierea de secvenţe de nucleotide este următoarea:

$$N = \begin{array}{c|ccccc} & A & T & C & G & - \\ \hline A & 1 & -1 & -1 & -1 & -5 \\ T & -1 & 1 & -1 & -1 & -5 \\ C & -1 & -1 & 1 & -1 & -5 \\ G & -1 & -1 & -1 & 1 & -5 \\ - & -5 & -5 & -5 & -5 & 0 \end{array}$$

Scorul unei alinieri de două secvenţe se poate calcula adunând valorile corespunzătoare din matrice pentru fiecare două nucleotide comparate.

Astfel, pentru următoarea aliniere:

```

ATTCAGG----ATTT
-T--CAGGGCCTATTT

```

scorul va fi:

$$S = N(A, -) + N(T, T) + N(T, -) + N(T, -) + N(C, C) + N(A, A) + N(G, G) + N(G, G) + N(-, G) + N(-, C) + N(-, C) + N(-, C) + N(-, T) + N(A, A) + N(T, T) + N(T, T)$$

$$T) + N(T, T) = -5 + 1 - 5 - 5 + 1 + 1 + 1 + 1 - 5 - 5 - 5 - 5 + 1 + 1 + 1 + 1 = 9 - 35 = -26$$

Un exemplu de matrice utilizată pentru alinierea de secvențe de aminoacizi este matricea PAM (Point Accepted Mutation) care a fost calculată observând diferențe dintre proteine înrudite. Alt exemplu este reprezentat de matricile BLOSUM (BLOCKS of amino acid SUBstitution Matrix) care se utilizează, de asemenea, pentru alinierea secvențelor de aminoacizi. Probabilitățile utilizate în calcularea matricilor BLOSUM sunt calculate pe baza blocurilor de secvențe conservate găsite în alinieri multiple de proteine. Cea mai utilizată matrice din seria BLOSUM este BLOSUM62.

Utilitare pentru alinierea de secvențe

1. *CLUSTAL-W* este un program care realizează alinieri globale între două sau mai multe secvențe de aminoacizi sau nucleotide. Se utilizează pentru compararea anumitor secvențe.
2. *BLAST* (Basic Local Alignment Search Tool) este un program care realizează tot alinieri, dar se utilizează pentru alinieri locale, fiind foarte util pentru căutarea de secvențe de ADN sau proteine în bazele de date. Astfel, se pot determina funcțiile anumitor secvențe pe baza similitudinii cu alte secvențe cunoscute și, de asemenea, se pot studia relații filogenetice între diferite specii.

2. Aplicații. CLUSTAL-W și BLAST

1. Să se acceseze *CLUSTAL-W* la următoarea adresă:

<http://www.ebi.ac.uk/Tools/clustalw2/index.html>

- a. Să se testeze gradul de asemănare între ribonucleazele pancreatice din fișierele FASTA pe baza scorului afișat de program.
- b. Cât este scorul pentru alinierea ribonucleazei de la om cu cea de la cimpanzeu? Dar de la om cu cea de la cangur?
- c. Să se calculeze scorurile obținute în urma unei alinieri multisevență (cel puțin 3 ribonucleaze diferite).

2. Să se acceseze următoarea adresă:

<http://www.uniprot.org/>

Să se caute secvențele în format FASTA pentru hemoglobină (subunitatea alpha) de la om și cea de la șoarece și să se compare utilizând CLUSTAL-W.

3. Să se caute pe *UniProt* secvența proteinei cu numărul de identificare *P26367*. Este o proteină codificată de gena *PAX-6* care controlează dezvoltarea ochilor la numeroase specii de organisme.

- a. Să se salveze secvența proteinei în format FASTA.
- b. Să se acceseze *PSI-BLAST* și să se introducă secvența proteinei:

<http://blast.ncbi.nlm.nih.gov/>

Programul va găsi toate organismele din baza de date la care gena *PAX-6* este prezentă. Căutarea se realizează în baza de date cu proteine și se afișează alinierea cu scorul cel mai bun.

3. Algoritmi pentru alinierea de secvențe

A. Un algoritm euristic pentru alinierea a două secvențe

Algoritmul de mai jos rezolvă problema alinierii globale a doua secvențe printr-o metodă “greedy” euristică. Un algoritm “greedy” urmărește rezolvarea problemei prin alegerea la fiecare pas a optimului local, sperând că în final se ajunge la optimul global.

Acest algoritm nu este simetric, uneori întoarce alinieri diferite în funcție de ordinea în care sunt date secvențele ca parametri. Apelul *align_asym(seq1, seq2)* poate întoarce alt rezultat decât apelul *align_asym(seq2, seq1)*.

```
def align_asym(seq1, seq2):
    i = 0
    while i < len(seq1) and i < len(seq2):
        if seq1[i] == seq2[i]:
            i += 1
        else:
            k = i
            min_idx1 = len(seq1)
            min_idx2 = len(seq2)
            while k < min_idx2 and k < len(seq1):
                idx = seq2.find(seq1[k], i)
                if idx >= 0 and min_idx2 > idx:
                    min_idx1 = k
                    min_idx2 = idx
                k += 1
            if min_idx1 > min_idx2:
                seq2 = seq2[:min_idx2] \
                    + '-' * (min_idx1 - min_idx2) \
                    + seq2[min_idx2:]
            elif min_idx1 < min_idx2:
                seq1 = seq1[:min_idx1] \
                    + '-' * (min_idx2 - min_idx1) \
                    + seq1[min_idx1:]

            i = max(min_idx1, min_idx2) + 1

    i = len(seq1)
    j = len(seq2)
    if i < j:
        seq1 = seq1 + '-' * (j - i)
    else:
        seq2 = seq2 + '-' * (i - j)

    return (seq1, seq2)
```

Se exemplifică funcționarea algoritmului pe două secvențe:

secv1: RQASPQT

secv2: RTSPTA

Algoritmul funcționează astfel:

i = contor de parcurgere secvențe

cât timp nu s-a ajuns la sfârșitul nici unei secvențe

dacă caracterele curente din cele două secvențe sunt egale, atunci:

$i += 1$

altfel:

$k \leftarrow i$;

$min_idx1 = \text{lungimea } secv1$;

$min_idx2 = \text{lungimea } secv2$;

se parcurge *sec1* de la *i* către sfârșit cu contorul *k*, fără ca valoarea *k* să depășească min_idx2

se caută caracterul de pe poziția *k* din *secv1* în *secv2*

dacă îl găsește și indexul din *secv2* (*idx*) este mai mic decât min_idx2 (cel salvat pentru perechea anterioară) atunci:

$min_idx1 \leftarrow k$;

$min_idx2 \leftarrow idx$;

(salvează indicii noii perechi)

se aliniază caracterele prin introducerea de caractere "-";

dacă una din secvențe este mai scurtă decât cealălaltă, la sfârșitul alinierii, se completează secvența mai scurtă cu caractere "-".

În urma alinierii celor două secvențe, se va obține următorul rezultat:

RQASPQT-

RT-SP-TA

Alte alinieri posibile mai proaste ar putea fi:

RQASPQT----

R-----TSPTA

RQ---ASPQT

RTSPTA----

Aplicații. Algoritm euristic pentru alinierea a două secvențe

1. Să se scrie o funcție care afișează un șir de caractere format din "+" și "-", care arată dacă acele caractere de la poziția i din două secvențe deja aliniate coincid sau nu. De exemplu:

```
RQASPQT-  
RT-SP-TA  
+--+--+--
```

2. Să se testeze algoritmul “*greedy*” de aliniere pe secvențele a două ribonucleaze pancreatice de la specii diferite, citite din fișiere în format FASTA. De exemplu, să se compare cea de la om cu cea de la cimpanzeu și cea de la om cu cea de la balenă. Ce se observă, care sunt mai asemănătoare?
3. Algoritmul nu găsește întotdeauna cea mai bună aliniere. Să se găsească un exemplu de două șiruri pentru care se observă în mod evident că alinierea obținută nu este cea mai bună. Să se explice de ce algoritmul “*greedy*” nu găsește alinierea cea mai bună.
4. Să se scrie o funcție, *simple_score(seq1, seq2)*, care calculează scorul unei alinieri globale, calculat ca procentul de caractere aliniate corect din lungimea secvențelor aliniate: $(nr_caractere_aliniat / lungime_secv) * 100$.
5. Să se scrie o funcție, *align(seq1, seq2, score)*, care:
 - primește ca parametri două secvențe nealinate și o funcție care calculează scorul;
 - apelează *align_asym(seq1, seq2)* și *align_asym(seq2, seq1)* și decide care aliniere este mai bună (are scor mai mare);
 - funcția întoarce cele două secvențe obținute în urma alinierii mai bune.
6. Să se exemplifice funcționarea corectă a funcțiilor implementate, pe secvențe sugestive din fișiere FASTA. Să se compare pe rând ribonucleaza pancreatică de la om cu cele de la alte specii, salvate în fișiere FASTA. Cum diferă scorurile? Ce se observă?