Laborator 7 – Bioinformatică

Alinierea de secvențe - continuare

A. Biblioteca NumPy

NumPy este o bibliotecă dezvoltată pentru manipularea de tablouri multidimensionale (arrays). Tablourile sunt obiecte mutabile. Tipurile de date conținute în tablouri pot fi: *int*, *float*, *complex*, *double*, *byte*, *long*, *string*, etc.

Să se testeze următoarele exemple. Ce afișează?

Observație. Ca orice bibliotecă, NumPy trebuie inclusă în scriptul Python care o utilizează:

```
from numpy import *
```

1.1.1 Funcții de bază pentru manipularea de tablouri

```
# crearea unui vector dintr-o lista
a = array([1,2,3,4,5])

print a
# crearea unui vector de lungime 4 cu elementele 0.0
b = zeros(4)

print b
# crearea unui vector de lungime 4 cu elementele 1.0
c = ones(4)

print c
# crearea unei matrici
d = array([[1,2],[3,4]])
```

```
print d

# matrice de 4x3 cu toate elementele 0.0

f = zeros((4,3))

print f

# matrice de 4x3 cu toate elementele nr intregi 0

g = zeros((4,3),int)

print g

# matrice de 4x3 cu toate elementele nr complexe

h = zeros((4,3), complex)

print h
```

1.1.2 Funcții matematice de bază pe tablouri

Exemplul 2

```
a = array([1,5,4.3])
b = array([-0.9,3,4])
print a + b

print a-b
print a*b

print a/b
```

```
a = array([[1,2,3,4],[5,6,7,8],[9,10,11,12]])
print a
```

```
print a[0]
 print a[-1]
 print a[0][0]
 print a[1][2]
 print a[:2]
 print a[:,1]
 print a.shape #dimensiunea unei matrici
 c = array([[1,2,3,4],[5,6,7,8],[3]) \#eroare, de ce?
Exemplul 4
 a = array([1,2,3,4])
 b = array([5,6,7,8])
 # inmultire pe componente
 print a*b
 # produs scalar - inmultire vector linie a cu vector colana b
 print dot(a,b)
 # inmultire un vector coloana a cu un vector linie b
 print outer (a,b)
Exemplul 5
 from numpy.random import *
 a = ranf((4,3))
 print a
```

```
b = random_integers(0,5,(4,3))
print b
```

```
b = random_integers(0,5,(4,3))
print sqrt(b)

print b**2

print exp(b)

print log(b)
```

1.1.3 Copieri, comparări și slicing

Exemplul 7

Variabilele conțin referințe către obiectele din memorie. Tablourile sunt tot obiecte. Ce afișează codul de mai jos și de ce?

```
b = array ([5.6, 4.2, 3.8])
c = b
c[0] = 0
print b
```

Exemplul 8

Pentru a copia locația de memorie se utilizează funcția copy.

```
from copy import *

b = array([5.6, 4.2, 3.8])

c=copy(b)

c[0]=0
```

```
print b
```

Copia se mai poate realiza prin adunarea lui 0 la vectorul b care creează alt obiect, la altă locație de memorie.

```
b = array([5.6, 4.2, 3.8])

c=b+0

c[0]=0

print b

Exemplul 9

a = ranf(5)

b = ranf(5)

print a

print b
```

Alte funcții de comparare sunt: less, less_equal, greater_equal, equal.

Exemplul 10

```
a = ranf(5) - 0.5
print a
print sign(a)
```

print greater(a,b)

print greater(a,b).astype(int)

```
a = [True, False, False,True]
b = [True, False, True,False]
print logical_and(a,b)
```

```
print logical_or(a,b)
```

```
• Conversii:
```

```
c = array([[1,2,3],[4,5,6],[7,8,9]])

print c

c[0,2] = 10.5

c[0][1] = 0.5

print c

c = c + 0.0

c[0,2] = 10.5

print c

c = c.astype(int)

print c
```

```
a = arange(15) #intoarce un array de la 0 la 14

print a

a = a/4

print a

a = arange(15)

a = a/4.

print a

b = [3,5,14,9,2]
```

```
print a[b]
```

```
• Sortare:
```

```
a = ranf(5)

print a

a.sort()

print a

b = a[::-1]

print b

a = ranf(5)

print a

print a.argsort() #indicii sevcventei sortate
```

```
a = ranf((3,2))

print a

print ravel(a)

print

print a.reshape((2,3))

print

print a

print

b = ranf((2,2))
```

```
print b

print

c = ranf((1,2))

print c

print

print concatenate((a,b,c))
```

• Permutarea random a unui tablou:

```
a = ranf(5)
print a
shuffle(a)
print a
```

Exemplul 17

• Funcții statistice:

```
a = ranf((4,3))

print a

# maximul din tot tabloul

print a.max()

# intoarce un vector cu maximele de pe fiecare coloana

print a.max(0)

# intoarce un vector cu maximele de pe fiecare linie

print a.max(1)
```

```
print a.min()

print a.min(0)

print a.min(1)

print a.sum()

print a.sum(0)

print a.sum(1)

print a.mean()

print a.mean(0)

print a.mean(1)
```

1.1.4 Conversii de la tablouri la alte tipuri

Exemplul 18

```
a = ranf(5)

print a.tolist()

a = ranf((4,3))

print a.tolist()
```

```
a = ranf(5)
print str(a)+'este string'
print map(str,a)
```

B. Altă modalitate de calculare a scorului în alinierea de secvențe:

Pentru calcularea scorului alinierilor secvențelor se iau în calcul și alți factori biologici precum tipurile de aminoacizi (polari, nepolari, ionizați). Astfel, când doi aminoacizi diferă, contează în calculul scorului final și dacă ambii sunt încărcați pozitiv, de exemplu, sau dacă unul este încărcat, iar celălalt nu. Astfel, scorul se poate calcula ținând cont și de probabilitatea înlocuirii unui aminoacid cu altul în cadrul aceleiași familii de proteine.

Să se implementeze o funcție, *blosum_score(seq1, seq2, gap = -8)*, care calculează scorul a doua secvențe aliniate astfel:

- când se întâlneşte "-", se adună la scor valoarea gap = -8;
- se adună ponderile perechilor de aminoacizi aliniați;
- în matricea *BLOSUM*, elementul (*i,j*) reprezintă valoarea care se adună la scor dacă se întalnește aliniată perechea *PBET[i] PBET[j]*, unde *PBET* este un vector care conține codurile aminoacizilor în urmatoarea ordine: PBET ='ARNDCQEGHILKMFPSTWYV'
- se observă că pe diagonala principală a matricii *BLOSUM* se găsesc valori pozitive deoarece corespund alinierii aceluiași aminoacid în ambele secvențe.
- se pot utiliza funcția *sir.index('caracter')* care întoarce poziția caracterului în șirul *sir*.
- să se utilizeze următoarea matrice *BLOSUM* de scoruri:

```
BLOSUM = array([
[5,-2,-1,-2,-1,-1,0,-2,-1,-2,-1,-1,-3,-1,1,0,-3,-2,0],
[-2, 7, -1, -2, -1, 1, 0, -3, 0, -4, -3, 3, -2, -3, -3, -1, -1, -3, -1, -3],
[-1,-1, 7, 2,-2, 0, 0, 0, 1,-3,-4,-0,-2,-4,-2,-1, 0,-4,-2,-3],
[-2,-2, 2, 8,-4, 0, 2,-1,-1,-4,-4,-1,-4,-5,-1, 0,-1,-5,-3,-4],
[-1, -4, -2, -4, 13, -3, -3, -3, -3, -2, -2, -2, -2, -4, -1, -1, -5, -3, -1],
[-1,-1, 0, 0,-3, 7, 2,-2, 1,-3,-2, 2, 0,-4,-1,-0,-1,-1,-1,-3],
[-1, 0, 0, 2, -3, 2, 6, -3, 0, -4, -3, 1, -2, -3, -1, -1, -1, -3, -2, -3],
[0,-3,0,-1,-3,-2,-3,8,-2,-4,-4,-2,-3,-4,-2,0,-2,-3,-3,-4],
[-2, 0, 1, -1, -3, 1, 0, -2, 10, -4, -3, 0, -1, -1, -2, -1, -2, -3, -1, 4],
[-1, -4, -3, -4, -2, -3, -4, -4, -4, 5, 2, -3, 2, 0, -3, -3, -1, -3, -1, 4],
[-2, -3, -4, -4, -2, -2, -3, -4, -3, 2, 5, -3, 3, 1, -4, -3, -1, -2, -1, 1],
[-1, 3, 0, -1, -3, 2, 1, -2, 0, -3, -3, 6, -2, -4, -1, 0, -1, -3, -2, -3],
[-1,-2,-2,-4,-2, 0,-2,-3,-1, 2, 3,-2, 7, 0,-3,-2,-1,-1, 0, 1],
[-3, -3, -4, -5, -2, -4, -3, -4, -1, 0, 1, -4, 0, 8, -4, -3, -2, 1, 4, -1],
[-1, -3, -2, -1, -4, -1, -1, -2, -2, -3, -4, -1, -3, -4, 10, -1, -1, -4, -3, -3],
[1,-1, 1, 0,-1, 0,-1, 0,-1,-3,-3, 0,-2,-3,-1, 5, 2,-4,-2,-2],
[0,-1,0,-1,-1,-1,-1,-2,-2,-1,-1,-1,-1,-2,-1,2,5,-3,-2,0],
[-3, -3, -4, -5, -5, -1, -3, -3, -3, -3, -2, -3, -1, 1, -4, -4, -3, 15, 2, -3],
[-2,-1,-2,-3,-3,-1,-2,-3, 2,-1,-1,-2, 0, 4,-3,-2,-2, 2, 8,-1],
[0,-3,-3,-4,-1,-3,-3,-4,-4, 4, 1,-3, 1,-1,-3,-2, 0,-3,-1, 5]])
```

C. Algoritm de programare dinamică pentru alinierea globală de secvențe. Algoritmul Needleman-Wunsch

Pentru proiectarea algoritmului de aliniere prin metoda programării dinamice este nevoie de o analiză mai atentă a problemei pentru a identifica modul de descompunere a acesteia.

Problema dată este următoarea: Se dau două secvențe, x și y, nu neapărat de aceeași dimensiune, peste alfabetul Σ și o funcție de scor s, $s:\Sigma\times\Sigma\to Z$. Se dorește determinarea celei mai bune alinieri globale, folosind modelul liniar de penalizare a golurilor (gaps). Modelul liniar presupune o penalizare constantă pentru fiecare gol.

Observații:

- Nu se acceptă în aliniere corespondența gap gap, deoarece aceasta poate fi ignorată.
- Funcția scor, s, se reprezintă ca o matrice.

Tehnica de programare dinamică se aplică pentru probleme care pot fi descompuse în subprobleme ale căror soluții se suprapun. Optimul global este obținut din soluțiile optime ale subproblemelor în care este descompusă problema inițială. Pentru a nu recalcula soluțiile subproblemelor, acestea sunt reținute într-o structură de date, cel mai adesea o matrice sau un vector.

Pentru a descompune problema alinierii globale, se va rezolva o problemă mai generală decât cea inițială și anume:

Se dau secvențele x, y, peste alfabetul Σ și o funcție de scor, s. Se dorește obținerea celei mai bune alinieri globale a subsecvenței formate din primele i caractere din secvența x cu primele j caractere din secvența y, pentru $\forall i \in \{1, \ldots, |x|\}$ și $\forall j \in \{1, \ldots, |y|\}$.

Se notează cu D(i,j) scorul celei mai bune alinieri dintre subsecvența primelor i caractere din x și subsecvența primelor j caractere din y.

Se notează cu g penalizarea unui gap, g < 0, pentru a descuraja introducerea gap-urilor.

Este evident ca D(0,0)=0.

De asemenea, $D(i,0) = i \cdot g$ și $D(0,j) = j \cdot g$.

D(i,0) reprezintă scorul alinierii primelor i caractere din x cu o subsecvență de 0 caractere din y. Este clar că secvența din x este aliniată cu o secvență de gap-uri.

In același mod se obține și D(0, j) pentru alinierea subsecvenței primelor j caractere din y cu gap-uri.

Formula recurentă pentru D(i, j) este următoarea:

$$D(i, j) = \max \begin{cases} D(i-1, j-1) + s(x_i, y_j) \\ D(i-1, j) + g \\ D(i, j-1) + g \end{cases}, i \in \{1, ..., |x|\}, j \in \{1, ..., |y|\}$$

$$(1)$$

D(0,0), D(i,0), D(0,j) reprezintă condițiile inițiale pentru formula recurentă. x_i reprezintă al i-lea element din secvența x și y_j reprezintă al j-lea element din secvența y. Formula recurentă (1) are următoarea semnificație:

- Prima ramura a recurenței corespunde cazului în care s-a realizat corespondența dintre caracterul x_i și caracterul y_j . Este nevoie să se cunoască scorul alinierii lui x_i cu y_j , adică $s(x_i, y_j)$ și scorul celei mai bune alinieri dintre primele i-l caractere din secvența x și primele j-l caractere din secvența y, adică D(i-l, j-l).
- A doua ramură a recurenței corespunde cazului în care s-a realizat corespondența dintre caracterul x_i și un gap. Este nevoie să se cunoască scorul celei mai bune alinieri dintre primele i-l caractere din secvența x și primele j caractere din secvența y, adică D(i-l, j), la care se adaugă penalizarea de gap, g.
- A treia ramură a recurenței corespunde cazului în care s-a realizat corespondența

dintre caracterul y_j și un gap. Este nevoie să se cunoască scorul celei mai bune alinieri dintre primele i caractere din secvența x și primele j-l caractere din secvența y, adică D(i, j-l), la care se adaugă penalizarea de gap, g.

$$\begin{vmatrix} x_1 \dots x_{i-1} & x_i & x_1 \dots x_{i-1} & x_i & x_1 \dots x_i & x_1 \dots x_i & y_1 \dots y_{j-1} & y_j & y_1 \dots y_j & y_$$

Scorul alinierii globale optime a secvențelor x și y se va găsi în D(|x|,|y|). Structura de date utilizată pentru implementarea acestei recurențe este o matrice de dimensiune $(|x|+1)\times(|y|+1)$. Matricea va fi initializată cu condițiile inițiale și apoi va fi completată pe linii de la elementul (I,I) până la elementul (|x|,|y|).

Observație:

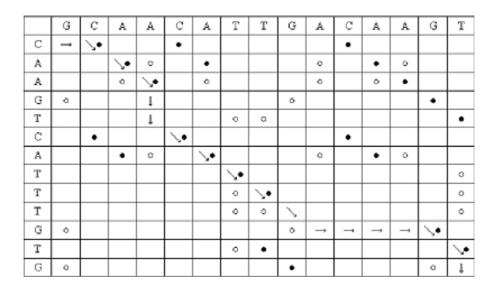
Anumite alinieri pot avea scoruri negative, datorită penalizarilor de gap sau a nepotrivirilor dintre caracterele aliniate.

Până acum s-a discutat despre obținerea scorului alinierii globale optime. Pentru obținerea alinierii propriu-zise, se pornește din elementul (|x|,|y|) din matricea D, și se reconstruiește alinierea dinspre ultimele caractere înspre primele, astfel:

- poziția inițială în matrice este (|x|, |y|);
- dacă valoarea elementului de la poziția curentă, (i, j), din matrice este egală cu rezultatul obținut din prima ramură a recurenței (1), atunci elementul curent va deveni (i-1, j-1), iar la aliniere se va adăuga la stânga corespondența dintre x_i și y_i .
- dacă valoarea elementului de la poziția curentă, (i, j), din matrice este egală cu rezultatul obținut din a doua ramură a recurenței (1), atunci elementul curent va deveni (i-1, j), iar la aliniere se va adăuga la stânga corespondența dintre x_i și un gap ("-").
- dacă valoarea elementului de la poziția curentă, (i, j), din matrice este egală cu rezultatul obținut din a treia ramură a recurenței (1), atunci elementul curent va deveni (i, j-1), iar la aliniere se va adăuga la stânga corespondența dintre un gap (" ") și y_j.
- se repetă algoritmul până se ajunge în elementul (0, 0).
- când se ajunge în elementul (0,0), s-a obținut alinierea corespunzătoare celui mai bun scor pentru alinierea globală a șirului x cu șirul y.

Observații:

- Dacă se ajunge într-o pozitie (i, 0), se adaugă, la stânga în aliniere, subsecvența primelor i caractere din x aliniate cu gap-uri.
- Dacă se ajunge într-o poziție (0, j), se adaugă, la stânga în aliniere, o secvență de *gap*-uri alinată cu subsecvența primelor j caractere din v.
- alinierea globală optimă a două secvențe se poate determina grafic astfel:



Aliniere globală optimă

D. Aplicații. Algoritmul Needleman-Wunsch

- 1. Să se completeze scheletul de cod din fisieru *alin_nw.py* si sa se testeze implementarea algoritmului Needleman-Wunsch.
- 2. În implementarea prezentată s-a utilizat o funcție de scor simplă, $simple(c_1, c_2)$, care întoarce valoarea 1 în caz de potrivire a două caractere și 0 în caz de nepotrivire. Să se scrie o funcție de scor, $BLOSUM(c_1, c_2)$, care să întoarcă valoarea corespunzatoare caracterelor c_1 și c_2 din matricea BLOSUM. Să se testeze implementarea algoritmului, utilizând această funcție de scor.
- 3. Sa se testeze algoritmul pe mai multe secvente ale proteinei ribonucleaza pancreatica si sa se compare rezultatele cu cele obtinute prin alinierea cu programul.