

▼ Obrada informacija

Laboratorijska vježba 2

U ovoj vježbi upoznat ćete se s jednom primjenom tehnika obrade informacija u bioinformatici. Ova laboratorijska vježba nosi 4 boda. Izvješće s ove laboratorijske vježbe potrebno je predati u .pdf formatu na *Moodle*. Izvješće koje predajete se mora zвати *Prezimelme.pdf*.

Osim biblioteka za rad s Fourierovom transformacijom (koristit ćemo samo numpy) koristit ćemo i biblioteku biopython koja sadrži puno korisnih alata iz područja bioinformatike. Mi ćemo je koristiti za jednostavnije baratanje bioinformatičkim tipovima podataka.

Biblioteka biopython dolazi s instalacijom Anaconde, ali ju je potrebno uključiti u okolinu (*environment*) koja se koristi.

Ako vježbu izvodite u Google Colab okruženju, morate instalirati biblioteku biopython. Instalaciju je potrebno izvršiti u sklopu prvog zadatka ove laboratorijske vježbe. Instalaciju izvodite sljedećim kodom:

```
try:
    import google.colab
    !pip install biopython
except ImportError:
    pass
```

Nakon izvođenja ovog koda, možete učitati biopython biblioteku.

1. Zadatak

Python biblioteke potrebne za laboratorijske vježbe su numpy i biopython. Uključite ih ("importirajte") i ispišite verziju svake od njih pomoću **[ime_biblioteke].version**.

UPUTA: Osnovna biopython biblioteka ima naziv Bio.

```
try:
    import google.colab
    !pip install biopython
except ImportError:
    pass
import Bio
import numpy as np

Collecting biopython
  Downloading biopython-1.81-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (3.1 MB)
    ━━━━━━━━━━━━━━━━ 3.1/3.1 MB 26.4 MB/s eta 0:00:00
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from biopython) (1.23.5)
Installing collected packages: biopython
Successfully installed biopython-1.81
```

Obrada informacija

Laboratorijska vježba 2

U ovoj vježbi upoznat ćete se s jednom primjenom tehnika obrade informacija u bioinformatici. Ova laboratorijska vježba nosi 4 boda. Izvješće s ove laboratorijske vježbe potrebno je predati u .pdf formatu na *Moodle*. Izvješće koje predajete se mora zвати *Prezimelme.pdf*.

Osim biblioteka za rad s Fourierovom transformacijom (koristit ćemo samo numpy) koristit ćemo i biblioteku biopython koja sadrži puno korisnih alata iz područja bioinformatike. Mi ćemo je koristiti za jednostavnije baratanje bioinformatičkim tipovima podataka.

Biblioteka biopython dolazi s instalacijom Anaconde, ali ju je potrebno uključiti u okolinu (*environment*) koja se koristi.

Ako vježbu izvodite u Google Colab okruženju, morate instalirati biblioteku biopython. Instalaciju je potrebno izvršiti u sklopu prvog zadatka ove laboratorijske vježbe. Instalaciju izvodite sljedećim kodom:

```
try:
    import google.colab
    !pip install biopython
except ImportError:
    pass
```

Nakon izvođenja ovog koda, možete učitati biopython biblioteku.

▼ Obrada informacija

Laboratorijska vježba 2

U ovoj vježbi upoznat ćete se s jednom primjenom tehnika obrade informacija u bioinformatici. Ova laboratorijska vježba nosi 4 boda. Izvješće s ove laboratorijske vježbe potrebno je predati u .pdf formatu na *Moodle*. Izvješće koje predajete se mora zвати *Prezimelme.pdf*.

Osim biblioteka za rad s Fourierovom transformacijom (koristit ćemo samo numpy) koristit ćemo i biblioteku biopython koja sadrži puno korisnih alata iz područja bioinformatike. Mi ćemo je koristiti za jednostavnije baratanje bioinformatičkim tipovima podataka.

Biblioteka biopython dolazi s instalacijom Anaconde, ali ju je potrebno uključiti u okolinu (*environment*) koja se koristi.

Ako vježbu izvodite u Google Colab okruženju, morate instalirati biblioteku biopython. Instalaciju je potrebno izvršiti u sklopu prvog zadatka ove laboratorijske vježbe. Instalaciju izvodite sljedećim kodom:

```
try:  
    import google.colab  
    !pip install biopython  
except ImportError:  
    pass
```

Nakon izvođenja ovog koda, možete učitati biopython biblioteku.

2. Zadatak

Uz laboratorijske vježbe dobili ste dvije datoteke s podacima. Datoteku koja sadrži referentni genom jednog soja bakterije *Escherichia coli* (*escherichia_coli_reference.fasta*) u FASTA formatu i datoteku koja sadrži skup očitanja dobivenih sekvenciranjem (*ecoli_ILL_small.fastq*) u FASTQ formatu.

Datoteke možete učitati koristeći metodu *parse()* iz biblioteke Bio.SeqIO. Metoda *parse()* vraća iterator koji možete pretvoriti u Python listu na sljedeći način:

```
reads = list(parse("ime_datoteke", "tip_datoteke"))
```

Tip datoteke postavite na "fasta" ili "fastq".

Učitajte obje datoteke te ispišite broj zapisa u svakoj od njih (broj elemenata u listi). Datoteka koja sadrži referencu trebala bi imati samo jedan zapis, dok bi datoteka s očitanjima trebala sadržavati veći broj zapisa.

NAPOMENA: Ako niste sigurni kako pronaći datoteke na disku iz Jupyter notebook-a, uvijek možete provjeriti radni direktorij sljedećim naredbama:

```
import os  
os.getcwd()
```

i promijeniti ga sa:

```
os.chdir()
```

Ako pak radite u Google Colab okruženju, koristite upute za učitavanje datoteka s Google diska iz prve laboratorijske vježbe.

```
from Bio.SeqIO import parse  
from google.colab import drive  
drive.mount('/content/drive')  
import os  
print(os.getcwd())  
os.chdir("drive/My Drive/Colab Notebooks/Obrada informacija/Lab2")  
genome_list = list(parse("escherichia_coli_reference.fasta", "fasta"))  
readings_list = list(parse("ecoli_ILL_small.fastq", "fastq"))  
print(len(genome_list))  
print(len(readings_list))
```

```
Mounted at /content/drive  
/content  
1  
38585
```

3. Zadatak

Svaki zapis koji ste učitali pomoću metode `Bio.SeqIO.parse()` sadrži Veći broj podataka od kojih su nam bitni samo neki. Naredbom print ispišite cijeli prvi zapis iz datoteke s očitanjima i iz datoteke s referencom.

Vidjet ćete da oba zapisa (među ostalim podacima) sadrže identifikator zapisa i sekvencu. Identifikator zapisa možete dohvatiti pomoću

```
zapis.id
```

dok sekvencu možete dohvatiti pomoću

```
zapis.seq
```

Ispišite identifikator i sekvencu za prvo očitanje te identifikator i prvih 200 znakova za referentni genom E.coli.

NAPOMENA: Referentni genom Escherichia coli je dugačak oko 4.5 milijuna slova

```
print(f"identifikator prvog očitanja {readings_list[0].id}")
reading_sequence = readings_list[0].seq
print(f"sekvenca prvog očitanja je {reading_sequence}")
print(f"identifikator E.coli {genome_list[0].id}")
reference_genome= genome_list[0].seq[:200:]
print(f"prvih 200 znakova za referentni E.coli je {reference_genome}")

identifikator prvog očitanja SRR2052522.671
sekvenca prvog očitanja je GATCTGGTGACCGGGTCGCGCAAAGTGATCATGCCATGGAACATTGCGCCAAAGATGGTTAGCAGAAAAATTGGGCCTCTGTATCATGCCACTCACTGCGCA/
identifikator E.coli NC_000913.3
prvih 200 znakova za referentni E.coli je AGCTTTTCAATTCTGACTGCAACGGGCAATATGTCCTGTGTTAAAAAAAGAGTGTCTGATAGCAGCTTCTGAACGGTTACCTGCC
```

4. Zadatak

Da bismo sekvence DNA analizirali metodama obrade signala, moramo pojednim nukleotidima (slovima) dodijeliti brojčane vrijednosti. Napišite funkciju u Pythonu koja će primiti slovo koje predstavlja nukleotid i vratiti odgovarajuću brojčanu vrijednost. Vrijednosti dodijelite na sljedeći način:

- A = 3
- G = 2
- C = -2
- T = -3

DNA sekvence mogu sadržavati i neke druge znakove (npr. 'N' koji označava da taj nuklotid nije poznat), njima dodijelite vrijednost 0. Također se može dogoditi da nukleotidi budu označeni i malim slovima, pa vodite računa da vaša funkcija mora vratiti ispravnu vrijednost i u tom slučaju.

```
def assign_number(letter):
    number = 0
    letter = letter.upper()
    if letter == "A":
        number = 3
    elif letter == "G":
        number = 2
    elif letter == "C":
        number = -2
    elif letter == "T":
        number = -3
    return number
```

5. Zadatak

Upotrebite napisanu funkciju da bi od prvog očitanja i od reference kreirali signal. Izračunajte korelaciju pomoću Fourierove transformacije. Zanemarite imaginarne vrijednosti.

```
nuc2num = ['A', 'G', 'C', 'T', 'N']
values = [assign_number(i) for i in nuc2num]
avg = np.average(values)
std = np.std(values)

x1 = [assign_number(i) for i in genome_list[0].seq]
x2 = [assign_number(i) for i in reading_sequence]
k_arr = range(-len(x2)+1, len(x1))
x1 = [(x-avg)/std for x in x1]
x2 = [(x-avg)/std for x in x2]
padding1 = [0]*(len(x2)-1)
padding2 = [0]*(len(x1)-1)
X1 = np.fft.fft(padding1 + x1)
```

```
X2 = np.fft.fft(x2 + padding2)
Cor = np.conjugate(X2)*X1
cor = np.fft.ifft(Cor)

k1 = k_arr[cor.argmax()]
print(f"Correlation by FFT: {np.real(cor)}")
print("k={}".format(k1))

Correlation by FFT: [-1.15384615  0.38461538 -0.19230769 ... -1.73076923 -2.30769231
-0.76923077]
k=2324486
```

6. Zadatak

Ispišite duljinu reference. Koristeći metode biblioteke *numpy*, izračunajte srednju vrijednost i standardnu devijaciju duljine očitanja (uzmite u obzir sva očitanja).

Primijetit ćete da su sva očitanja jednake duljine.

```
print(len(genome_list[0].seq))
avg_list = []
for i in range(len(readings_list)):
    avg_list.append(len(readings_list[i].seq))
print(np.average(avg_list))
print(np.std(avg_list))

4641652
121.0
0.0
```

7. zadatak

Što ako želimo izračunati korelaciju za veći broj očitanja i istu referencu? To je tipičan slučaj u bioinformatici jer uređaji za sekvenciranje proizvode tisuće i desetke tisuća očitanja koja se potom mapiraju na istu referencu.

Ako korelaciju računamo izravno, potrebno ju je svaki put izračunati iz početka. Ako korelaciju računamo pomoću FFT-a, transformaciju za referencu potrebno je napraviti samo jednom.

Izračunajte korelacije za prvih 10 očitanja.

```
nuc2num = ['A', 'G', 'C', 'T', 'N']
values = [assign_number(i) for i in nuc2num]
avg = np.average(values)
std = np.std(values)
x1 = [assign_number(i) for i in genome_list[0].seq]
x1 = [(x-avg)/std for x in x1]
padding1 = [0] * (len(readings_list[0].seq) - 1)
padding2 = [0]*(len(x1)-1)
X1 = np.fft.fft(padding1+x1)

for i in range(10):

    sequence = readings_list[i].seq
    x2 = [assign_number(i) for i in sequence]
    x2 = [(x-avg)/std for x in x2]
    X2 = np.fft.fft(x2 + padding2)
    #print(len(X2))
    #print(len(X1))
    Cor = np.conjugate(X2)*X1
    cor = np.fft.ifft(Cor)
    k = k_arr[cor.argmax()]
    print(f"Correlation by FFT: {np.real(cor)}")
    print("k={}".format(k))

    Correlation by FFT: [-1.15384615  0.38461538 -0.19230769 ... -1.73076923 -2.30769231
-0.76923077]
    k=2324486
    Correlation by FFT: [ 1.73076923e+00 -3.55272330e-15 -3.65384615e+00 ... -1.73076923e+00
-2.5000000e+00 -1.15384615e+00]
    k=1877260
    Correlation by FFT: [-1.15384615 -2.5           -2.11538462 ... -3.46153846 -1.92307692
-0.76923077]
    k=557777
    Correlation by FFT: [ 1.73076923 -0.57692308 -1.15384615 ... -1.73076923 -0.96153846
-1.15384615]
    k=1144877
    Correlation by FFT: [ 1.73076923e+00 -1.59872122e-14 -3.07692308e+00 ...  1.92307692e-01
-9.61538462e-01 -1.15384615e+00]
    k=3583639
```

```
Correlation by FFT: [ 1.15384615  1.92307692  1.73076923 ... -4.03846154 -2.5
-1.15384615]
k=4051518
Correlation by FFT: [-1.15384615e+00  9.61538462e-01  3.07692308e+00 ... -1.73076923e+00
 3.55272209e-15  7.69230769e-01]
k=2293706
Correlation by FFT: [ 1.15384615  2.5       -1.34615385 ...  1.73076923  0.96153846
 1.15384615]
k=1011323
Correlation by FFT: [-1.73076923 -2.88461538 -1.73076923 ... -2.30769231 -2.88461538
-1.15384615]
k=628546
Correlation by FFT: [-1.15384615 -1.92307692  1.15384615 ... -1.15384615 -0.57692308
-1.15384615]
k=1497921
```

8. zadatak

Na temelju najveće vrijednosti korelacije između reference i prvog očitanja pronađite poziciju na referenci koja je najsličnija očitanju. Pozicija odgovara vrijednosti parametra k za koji je korelacija najveća.

Napišite metodu koja će primiti dva niza znakova jednake duljine, usporediti znakove na istim pozicijama i vratiti broj razlika (Hammingova udaljenost).

"Izrežite" dio reference koji je najsličniji očitanju (iste duljine kao i očitanje) i usporedite ga s očitanjem pomoću napisane funkcije.

```
def Hamming_distance(sequence1, sequence2):
    distance = 0
    for i in range(0,len(sequence1)):
        if sequence2[i] != sequence1[i]:
            distance += 1
    return distance

part_ref = genome_list[0].seq[k1:k1+len(readings_list[0].seq)]
print(f"razlika u Hammingovoj distanci je {Hamming_distance(part_ref, readings_list[0].seq)}")

razlika u Hammingovoj distanci je 9
```

9. zadatak

U datoteci "ecoli_ILL_small_aln.sam" dana su već izračunata poravnanja svih očitanja na referencu u SAM formatu. SAM je tekstualni "tab separated" format. U prvom stupcu se nalazi identifikator očitanja, dok se u četvrtom stupcu nalazi pozicija na referenci na koju je očitanje najbolje poravnato (ostali stupci ne zanimaju). Također, datoteka s poravnanjima sadrži i nekoliko header readaka kojima prvi stupac počinje sa znakom '@', njih također možete zanemariti.

Otvorite datoteku s poravnanjima i pronađite poravnanje za prvo očitanje (identifikator očitanja u datoteci s očitanjima i datoteci s poravnanjima mora biti jednak). Usporedite poziciju u datoteci sa pozicijom koju ste dobili pomoću korelacije.

UPOUTA: TSV datoteke možete otvoriti na sljedeći način:

```
tsv_file = open("file_name")
tsv_rows = csv.reader(tsv_file, delimiter="\t")
```

Varijabla tsv_rows će sadržavati listu redaka, a svaki redak biti lista vrijednosti (po jedna za svaki stupac).

```
import csv
tsv_file = open("ecoli_ILL_small_aln.sam")
tsv_rows = csv.reader(tsv_file, delimiter="\t")
tsv_list_rows = list(tsv_rows)
for i in tsv_list_rows:
    if i[0] == readings_list[0].id:
        print(f"vrijednost u SAM datoteci je {i[3]}")
print(f"korelacijom sam dobio {k1}")

vrijednost u SAM datoteci je 2324487
korelacijom sam dobio 2324486
```

10. zadatak

Za prvo očitanje pozicija dobivena pomoću korelacije trebala bi biti 2324486, dok je pozicija u datoteci s poravnanjima 2324487. Razlikuju se samo za 1 pa možemo zaključiti da nam je korelacija dala dobru poziciju za poravnanje.

Prisjetimo se da korelacija ne računa točno poravnanje već ju koristimo samo da bi našli kandidatne pozicije za točno računanje. Tek onda na takvim pozicijama možemo točno poravnanje izračunati pomoću algoritma dinamičkog programiranja. Ako bi primijenili dinamičko

programiranje za računanje poravnajanja očitanja s cijelom referencom, postupak bi bio znatno sporiji i zahtijevao bi veliku količinu radne memorije.

Ako želite to možete isprobati pomoću algoritama za poravnanje u biblioteci *bioparser*. Lokalno poravnanje možete izračunati metodom:

```
Bio.pairwise2.align.localxx(seq1, seq2)
```

Za prvi 100 očitanja izračunajte korelaciju te pomoću korelacije poziciju najveće sličnosti očitanja i reference. Usporedite rezultat sa podacima u datoteci s poravnanjima. Ispišite broj očitanja za koja se dvije pozicije razlikuju za najviše 5 mesta.

```
nuc2num = ['A', 'G', 'C', 'T', 'N']
values = [assign_number(i) for i in nuc2num]
avg = np.average(values)
std = np.std(values)
x1 = [assign_number(i) for i in genome_list[0].seq]
x1 = [(x-avg)/std for x in x1]
padding1 = [0] * (len(readings_list[0].seq) - 1)
padding2 = [0]*(len(x1)-1)
X1 = np.fft.fft(padding1+x1)
number_of_readings = 0
for i in range(100):
    sequence = readings_list[i].seq
    x2 = [assign_number(i) for i in sequence]
    x2 = [(x-avg)/std for x in x2]
    X2 = np.fft.fft(x2 + padding2)
    Cor = np.conjugate(X2)*X1
    cor = np.fft.ifft(Cor)
    k = k_arr[cor.argmax()]
    for row in tsv_list_rows:
        if row[0] == readings_list[i].id:
            if abs(int(row[3]) - k) <= 5:
                number_of_readings+=1
print(f"Broj očitanja za koja se dvije pozicije razlikuju za najviše 5 mesta je {number_of_readings}")
```

Broj očitanja za koja se dvije pozicije razlikuju za najviše 5 mesta je 52

11. ZAKLJUČAK

Očekivani broj točno pozicioniranih očitanja je 50, jer smo za sada uspješno radili samo s očitanjima na jednom lancu DNA.

Prolaskom kroz zadatke u ovoj vježbi dobili ste kratak uvod u rad s bioinformatičkim podacima i tehnikama obrade signala.