

2LabVježba_OI

November 16, 2023

1 Obrada informacija

1.1 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 zvati *PrezimeIme.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.

```
[1]: # Ovo je mjesto na kojem možete izvoditi svoj kod.
import numpy as np
import Bio
from Bio import SeqIO

print(f"numpy {np.__version__}")
print(f"biopython {Bio.__version__}")
```

numpy 1.24.3
biopython 1.78

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 očitavanja 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.

```
[2]: # Ovo je mjesto na kojem možete izvoditi svoj kod.
file_path_fasta = "/Users/dominik/Desktop/Obrada informacija/Lab2/
↳escherichia_coli_reference.fasta"

file_path_fastq = "/Users/dominik/Desktop/Obrada informacija/Lab2/
↳ecoli_ILL_small.fastq"

reference = list(SeqIO.parse(file_path_fasta, "fasta"))
reads = list(SeqIO.parse(file_path_fastq, "fastq"))

print(f"Broj zapisa datoteke koja sadrži referencu: {len(reference)}")
print(f"Broj zapisa datoteke koja sadrži očitavanja: {len(reads)}")
```

Broj zapisa datoteke koja sadrži referencu: 1
Broj zapisa datoteke koja sadrži očitavanja: 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 očitavanje te identifikator i prvih 200 znakova za referentni genom E.coli.

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

```
[3]: # Ovo je mjesto na kojem možete izvoditi svoj kod.
print(f"Prvi zapis datoteke s očitanjima\n{reads[0]}\n")
print(f"Prvi(jedini) zapis datoteke s referencom\n{reference[0]}\n")

print(f"Prvo očitanje\nID: {reads[0].id}\nSeq: {reads[0].seq}\n")
print(f"Referentni genom prvih 200\nID: {reference[0].id}\nSeq: {reference[0].
↪seq[:200]}\n")
```

Prvi zapis datoteke s očitanjima

ID: SRR2052522.671

Name: SRR2052522.671

Description: SRR2052522.671 HWI-EAS390_0001:4:1:6915:1123/1

Number of features: 0

Per letter annotation for: phred_quality

Seq('GATCTGGTGACCGGGTCGCGCAAAGTGATCATCGCCATGGAACATTGCGCCAAA...TGC')

Prvi(jedini) zapis datoteke s referencom

ID: NC_000913.3

Name: NC_000913.3

Description: NC_000913.3 Escherichia coli str. K-12 substr. MG1655, complete genome

Number of features: 0

Seq('AGCTTTTCATTCTGACTGCAACGGGCAATATGTCTCTGTGTGGATTAAAAAAAG...TTC')

Prvo očitanje

ID: SRR2052522.671

Seq: GATCTGGTGACCGGGTCGCGCAAAGTGATCATCGCCATGGAACATTGCGCCAAAGATGGTTCAGCAAAATTTTG
GGCCTCTGTATCATGCCACTCACTGCGCAATATCCGGATCAAATGC

Referentni genom prvih 200

ID: NC_000913.3

Seq: AGCTTTTCATTCTGACTGCAACGGGCAATATGTCTCTGTGTGGATTAAAAAAAGAGTGTCTGATAGCAGCTTCTG
AACTGGTTACCTGCCGTGAGTAAATTAATTTTATTGACTTAGGTCATAAATACTTTAACCAATATAGGCATAGCGCA
CAGACAGATAAAAATTACAGAGTACACAACATCCATGAAACGCAT

4. Zadatak

Da bismo sekvence DNA analizirali metodama obrade signala, moramo pojedinim 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 nukleotid 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.

[4]: *# Ovo je mjesto na kojem možete izvoditi svoj kod.*

```
nuc2num = {}
nuc2num['A'] = 3
nuc2num['G'] = 2
nuc2num['C'] = -2
nuc2num['T'] = -3

def parse_nuc(nuc):
    value = 0
    if nuc.upper() in nuc2num:
        value = nuc2num.get(nuc.upper())
    return value

# testing

# seq1='AGCTNagctNMOP'

# string = ''
# for s in seq1:
#     val = parse_nuc(s)
#     string = string + str(val) + ' '

# print(string)

# output is 3 2 -2 -3 0 3 2 -2 -3 0 0 0 0
```

5. Zadatak

Upotrebite napisanu funkciju da bi od prvog očitavanja i od reference kreirali signal. Izračunajte korelaciju pomoću Fourierove transformacije. Zanimajte imaginarnu vrijednost.

[5]: *# Ovo je mjesto na kojem možete izvoditi svoj kod.*

```
signal_reference = [parse_nuc(nuc) for nuc in reference[0].seq]
signal_first = [parse_nuc(nuc) for nuc in reads[0].seq]
```

```

len_reference = len(signal_reference)
len_first = len(signal_first)
k_arr = range(-len_first + 1, len_reference)

avg = np.average(list(nuc2num.values()))
std = np.std(list(nuc2num.values()))

signal_reference = [(x-avg)/std for x in signal_reference]
signal_first = [(x-avg)/std for x in signal_first]

padding_reference = [0] * (len_first - 1)
padding_first = [0] * (len_reference - 1)

X_ref = np.fft.fft(padding_reference + signal_reference)
X_fst = np.fft.fft(signal_first + padding_first)
Cor = np.conjugate(X_fst) * X_ref
cor = np.fft.ifft(Cor)

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

```

Correlation by FFT:

```

[-0.92307692-4.44089210e-15j  0.30769231+1.86213181e-15j
 -0.15384615-5.05601992e-15j ... -1.38461538+3.57227567e-15j
 -1.84615385-8.04796621e-15j -0.61538462-1.03855894e-14j]
k = 2324486

```

6. Zadatak

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

Primijetit ćete da su sva očitavanja jednake duljine.

```

[6]: # Ovo je mjesto na kojem možete izvoditi svoj kod.
print("Duljina reference: {}".format(len_reference))

len_reads = [len(x.seq) for x in reads]

avg_reads = np.mean(len_reads)
std_reads = np.std(len_reads)

print(f"Srednja vrijednost duljine očitavanja: {avg_reads}")
print(f"Standardna devijacija duljine očitavanja: {std_reads}")

```

Duljina reference: 4641652

Srednja vrijednost duljine očitavanja: 121.0

Standardna devijacija duljine očitavanja: 0.0

7. zadatak

Što ako želimo izračunati korelaciju za veći broj očitavanja i istu referencu? To je tipičan slučaj u bioinformatici jer uređaji za sekvenciranje proizvode tisuće i desetke tisuća očitavanja 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 očitavanja.

```
[7]: # Ovo je mjesto na kojem možete izvoditi svoj kod.
for i, read in enumerate(reads[:10]):
    signal_read = [parse_nuc(nuc) for nuc in read.seq]

    X_read = np.fft.fft(signal_read + padding_first)
    Cor = np.conjugate(X_read) * X_ref
    cor = np.fft.ifft(Cor)

    k = k_arr[cor.argmax()]
    print(f"Correlation by FFT reading {i + 1}:")
    print(cor)
    print("k = {}".format(k))
```

Correlation by FFT reading 1:

```
[-2.35339362+5.32907052e-15j  0.78446454-2.27121682e-14j
 -0.39223227-7.42103756e-15j ... -3.53009043+8.25889496e-15j
 -4.70678724+1.13112010e-14j -1.56892908-4.31368224e-15j]
```

k = 2324486

Correlation by FFT reading 2:

```
[ 3.53009043e+00+1.95399252e-14j  2.13162809e-14+1.77637127e-15j
 -7.45241314e+00-1.77847077e-15j ... -3.53009043e+00+2.06933144e-14j
 -5.09901951e+00-1.47615161e-14j -2.35339362e+00-1.58552422e-15j]
```

k = 1877260

Correlation by FFT reading 3:

```
[-2.35339362-7.10542736e-15j -5.09901951+6.35546990e-16j
 -4.31455497+1.82883080e-14j ... -7.06018086+1.62953504e-14j
 -3.9223227 +6.61348552e-15j -1.56892908+5.45630742e-15j]
```

k = 557777

Correlation by FFT reading 4:

```
[ 3.53009043-2.22044605e-14j -1.17669681-1.96665613e-14j
 -2.35339362+2.56396342e-14j ... -3.53009043+2.82426497e-14j
 -1.96116135+8.63578923e-15j -2.35339362+3.50518536e-14j]
```

k = 1144877

Correlation by FFT reading 5:

```
[ 3.53009043e+00+1.42108547e-14j  2.13163169e-14-5.15143339e-14j
 -6.27571632e+00+5.60785366e-15j ...  3.92232270e-01-7.62859843e-15j
 -1.96116135e+00+1.04121220e-14j -2.35339362e+00-1.97931809e-14j]
```

```

k = 3583639
Correlation by FFT reading 6:
[ 2.35339362+2.04281037e-14j  3.9223227 +1.61132484e-14j
  3.53009043-2.06958612e-14j ... -8.23687768+7.42799480e-15j
 -5.09901951-3.91864418e-14j -2.35339362-6.47052687e-15j]
k = 4051518
Correlation by FFT reading 7:
[-2.35339362e+00+2.13162821e-14j  1.96116135e+00-6.59832771e-15j
  6.27571632e+00+5.93858121e-15j ... -3.53009043e+00+1.15359474e-15j
  5.32909216e-15-7.99359135e-15j  1.56892908e+00+3.80644357e-14j]
k = 2293706
Correlation by FFT reading 8:
[ 2.35339362-1.06581410e-14j  5.09901951+3.80532467e-15j
 -2.74562589+1.42217456e-14j ...  3.53009043+1.83868120e-14j
  1.96116135-2.68435430e-14j  2.35339362-1.12930557e-14j]
k = 1011323
Correlation by FFT reading 9:
[-3.53009043+1.06581410e-14j -5.88348405-2.46139511e-14j
 -3.53009043-3.79239181e-14j ... -4.70678724+2.72209484e-15j
 -5.88348405-2.24983280e-14j -2.35339362-2.15695483e-14j]
k = 628546
Correlation by FFT reading 10:
[-2.35339362-1.33226763e-14j -3.9223227 -2.79058079e-15j
  2.35339362+3.20919343e-14j ... -2.35339362+1.20191668e-14j
 -1.17669681+4.82618366e-15j -2.35339362+9.51673067e-15j]
k = 1497921

```

8. zadatak

Na temelju najveće vrijednosti korelacije između reference i prvog očitavanja pronađite poziciju na referenci koja je najslbližnja očitavanju. 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 najslbližiji očitavanju (iste duljine kao i očitanje) i usporedite ga s očitanjem pomoću napisane funkcije.

```

[8]: # Ovo je mjesto na kojem možete izvoditi svoj kod.

def get_hamming_dist(seq1, seq2):
    return sum(seq1[i] != seq2[i] for i in range(len(seq1))) # vraca zbroj
    ↪ razlika između 2 stringa jednakih duljina

pos = 2324486 # ovo je korelacija između reference i prvog očitavanja

seq_cut = reference[0].seq[pos:pos + len(reads[0].seq)]

```

```
diff = get_hamming_dist(seq_cut, reads[0].seq)
print(f"Razlika između izrezanog dijela reference i prvog očitavanja: {diff}")
```

Razlika između izrezanog dijela reference i prvog očitavanja: 9

9. zadatak

U datoteci "ecoli_ILL_small_aln.sam" dana su već izračunata poravnanja svih očitavanja na referencu u SAM formatu. SAM je tekstualni "tab separated" format. U prvom stupcu se nalazi identifikator očitavanja, dok se u četvrtom stupcu nalazi pozicija na referenci na koju je očitavanje najbolje poravnato (ostali stupci nas 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 očitavanje (identifikator očitavanja u datoteci s očitavanjima 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).

```
[9]: # Ovo je mjesto na kojem možete izvoditi svoj kod.
import csv # potrebno importirati za citanje tsv filea
tsv_file = open("/Users/dominik/Desktop/Obrada informacija/Lab2/
↳ecoli_ILL_small_aln.sam")
tsv_rows = list(csv.reader(tsv_file, delimiter="\t"))
for line in tsv_rows:
    if line[0] == reads[0].id:
        print(f"Pozicija poravnanja iz datoteke = {line[3]}\nPozicija_
↳poravnanja iz korelacije = {pos}")
```

Pozicija poravnanja iz datoteke = 2324487

Pozicija poravnanja iz korelacije = 2324486

10. zadatak

Za prvo očitavanje 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 algoritama dinamičkog programiranja. Ako bi primijenili dinamičko programiranje za računanje poravnanja očitavanja 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 prvih 100 očitavanja izračunajte korelaciju te pomoću korelacije poziciju najveće sličnosti očitavanja i reference. Usporedite rezultat sa podacima u datoteci s poravnanjima. Ispišite broj očitavanja za koja se dvije pozicije razlikuju za najviše 5 mjesta.

```
[10]: # Ovo je mjesto na kojem možete izvoditi svoj kod.

# from Bio import pairwise2
# alignments = pairwise2.align.localxx('CAGGTAC', 'AGGTAC')
# print(alignments)
# [Alignment(seqA='CAGGTAC', seqB='-AGGTAC', score=6.0, start=1, end=7)]

counter = 0
for i, read in enumerate(reads[:100]):
    signal_read = [parse_nuc(nuc) for nuc in read.seq]

    X_read = np.fft.fft(signal_read + padding_first)
    Cor = np.conjugate(X_read) * X_ref
    cor = np.fft.ifft(Cor)

    k = k_arr[cor.argmax()]

    for line in tsv_rows:
        if line[0] == read.id:
            if np.abs(int(line[3]) - k) < 6:
                counter = counter + 1

print(f"Broj očitavanja za koja se 2 pozicije razlikuju za najviše 5 mjesta:␣
↪{counter}")
```

Broj očitavanja za koja se 2 pozicije razlikuju za najviše 5 mjesta: 50

11. ZAKLJUČAK

Očekivani broj točno pozicioniranih očitavanja 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.