

HaralovicMarko

November 10, 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.

```
[ ]: import Bio
import numpy as np
print(Bio.__version__)
print(np.__version__)
```

1.81

1.23.5

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.

```
[ ]: from Bio.SeqIO import parse

referentni_genom = list(parse("./escherichia_coli_reference.fasta", "fasta"))
ocitanja_sekvenciranja = list(parse("./ecoli_ILL_small.fastq", "fastq"))

print(f"Duljina referentnog genoma je {len(referentni_genom)}")
print(f"Duljina očitane sekvence je {len(ocitanja_sekvenciranja)}")
```

```
Duljina referentnog genoma je 1
Duljina očitane sekvence je 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

```
[ ]: print(f"Identifikator zapisa : {referentni_genom[0].id}")
      print(f"Sekvenca zapisa: {referentni_genom[0].seq[:200]}")
      print("-----")
      print(f"ID sekvence : {ocitanja_sekvenciranja[0].id}")
      print(f"Sekvenca prvog očitavanja: {ocitanja_sekvenciranja[0].seq}")
```

Identifikator zapisa : NC_000913.3

Sekvenca zapisa: AGCTTTTCATTCTGACTGCAACGGGCAATATGTCTCTGTGTGGATTAAAAAAGAGTGTCTGA
TAGCAGCTTCTGAACTGGTTACCTGCCGTGAGTAAATTTAAATTTTATTGACTTAGGTCACTAAATACTTTAACCAATAT
AGGCATAGCGCACAGACAGATAAAAATTACAGAGTACACAACATCCATGAAACGCAT

ID sekvence : SRR2052522.671

Sekvenca prvog očitavanja: GATCTGGTGACCGGGTCGCGCAAAGTGATCATCGCCATGGAACATTGCGCCAAAG
ATGGTTCAGCAAAAATTTTGGGCCTCTGTATCATGCCACTCACTGCGCAATATCCGGATCAAATGC

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.

```
[ ]: nucleotide_letter_to_digit = {
      'A':3,
      'G':2,
      'C':-2,
      'T':-3
    }

def enumerate_nucleotide(nucleotide):
    nucleotide = nucleotide.upper()
    try:
        number = nucleotide_letter_to_digit[nucleotide]
    except KeyError:
        number = 0
    return number
```

5. Zadatak

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

```
[ ]: referentni_genom_signal = [enumerate_nucleotide(nucleotide) for nucleotide in
    ↪referentni_genom[0].seq]
ocitanja_sekvenciranja_signal = [enumerate_nucleotide(nucleotide) for nucleotide
    ↪in ocitanja_sekvenciranja[0].seq]
```

```
[ ]: print(f"Signal genoma: {referentni_genom_signal[:200]}")
print("-----")
print(f"Signal očitavanja: {ocitanja_sekvenciranja_signal}")
```

```
Signal genoma: [3, 2, -2, -3, -3, -3, -3, -2, 3, -3, -3, -2, -3, 2, 3, -2, -3,
2, -2, 3, 3, -2, 2, 2, 2, -2, 3, 3, -3, 3, -3, 2, -3, -2, -3, -2, -3, 2, -3, 2,
-3, 2, 2, 3, -3, -3, 3, 3, 3, 3, 3, 3, 2, 3, 2, -3, 2, -3, -2, -3, 2, 3, -3,
3, 2, -2, 3, 2, -2, -3, -3, -2, -3, 2, 3, 3, -2, -3, 2, 2, -3, -3, 3, -2, -2,
-3, 2, -2, -2, 2, -3, 2, 3, 2, -3, 3, 3, 3, -3, -3, 3, 3, 3, 3, -3, -3, -3, -3,
3, -3, -3, 2, 3, -2, -3, -3, 3, 2, 2, -3, -2, 3, -2, -3, 3, 3, 3, -3, 3, -2, -3,
-3, -3, 3, 3, -2, -2, 3, 3, -3, 3, -3, 3, 2, 2, -2, 3, -3, 3, 2, -2, 2, -2, 3,
-2, 3, 2, 3, -2, 3, 2, 3, -3, 3, 3, 3, 3, 3, -3, -3, 3, -2, 3, 2, 3, 2, -3, 3,
-2, 3, -2, 3, 3, -2, 3, -3, -2, -2, 3, -3, 2, 3, 3, 3, -2, 2, -2, 3, -3]
```

```
-----
Signal očitavanja: [2, 3, -3, -2, -3, 2, 2, -3, 2, 3, -2, -2, 2, 2, 2, -3, -2, 2,
-2, 2, -2, 3, 3, 3, 2, -3, 2, 3, -3, -2, 3, -3, -2, 2, -2, -2, 3, -3, 2, 2, 3,
3, -2, 3, -3, -3, 2, -2, 2, -2, -2, 3, 3, 3, 2, 3, -3, 2, 2, -3, -3, -2, 3, 2,
-2, 3, 3, 3, 3, 3, -3, -3, -3, -3, 2, 2, 2, -2, -2, -3, -2, -3, 2, -3, 3, -3,
-2, 3, -3, 2, -2, -2, 3, -2, -3, -2, 3, -2, -3, 2, -2, 2, -2, 3, 3, -3, 3, -3,
-2, -2, 2, 2, 3, -3, -2, 3, 3, 3, -3, 2, -2]
```

```
[ ]: avg_ref = np.average(referentni_genom_signal)
std_ref = np.std(referentni_genom_signal)
referentni_genom_signal = [(x - avg_ref) / std_ref for x in
    ↪referentni_genom_signal]

avg_ocitanja = np.average(ocitanja_sekvenciranja_signal)
std_ocitanja = np.std(ocitanja_sekvenciranja_signal)
ocitanja_sekvenciranja_signal = [(x - avg_ocitanja) / std_ocitanja for x in
    ↪ocitanja_sekvenciranja_signal]
```

```
[ ]: padding_ref = [0] * (len(ocitanja_sekvenciranja_signal) - 1)
padding_ocitanja = [0] * (len(referentni_genom_signal) - 1)
k_arr =
    ↪range(-len(ocitanja_sekvenciranja_signal)+1, len(referentni_genom_signal))
k_arr
```

```
[ ]: range(-120, 4641652)
```

```
[ ]: FFT_ref = np.fft.fft(padding_ref + referentni_genom_signal)
      FFT_ocitanja = np.fft.fft(ocitanja_sekvenciranja_signal + padding_ocitanja)
```

```
[ ]: cross_correlation_freq = np.conjugate(FFT_ocitanja) * FFT_ref
      cross_correlation = np.fft.ifft(cross_correlation_freq)
      cross_correlation_magnitude = np.abs(cross_correlation)
      k = k_arr[cross_correlation.argmax()]
```

```
[ ]: print("Cross-correlation FFT:")
      print(cross_correlation_magnitude)
      print("Indeks najveće korelacije: {}".format(k))
```

```
Cross-correlation FFT:
[0.96628788 0.24537804 0.19311781 ... 1.28947659 1.7916225 0.59293484]
Indeks najveće korelacije: 2324486
```

6. Zadatak

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

Primijetiti ćete da su sva očitavanja jednake duljine.

```
[ ]: print("Duljina genoma : " , len(referentni_genom[0].seq))

      duljina_ocitanja = [len(ocitanje.seq) for ocitanje in ocitanja_sekvenciranja]

      mean_length = np.mean(duljina_ocitanja)
      std_dev_length = np.std(duljina_ocitanja)

      print("Srednja duljina očitavanja :", mean_length)
      print("Srednja devijacija duljine očitavanja:", std_dev_length)
```

```
Duljina genoma : 4641652
Srednja duljina očitavanja : 121.0
Srednja devijacija duljine očitavanja: 0.0
```

7. zadatak

Što ako želimo izračunati korelaciju za veći broj očitavanja i istu referencu? To je tipičan slučaj u bioinformatičari 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.

```
[ ]: for i in range(10):
      ocitanje_signal = [enumerate_nucleotide(nucleotide) for nucleotide in
      ↪ ocitanja_sekvenciranja[i].seq]
      avg_ocitanja = np.average(ocitanje_signal)
```

```

std_ocitanja = np.std(ocitanje_signal)
ocitanje_signal = [(x - avg_ocitanja) / std_ocitanja for x in
↪ocitanje_signal]

padding_ocitanja = [0] * (len(referentni_genom_signal) - 1)
padding_ref = [0] * (len(ocitanje_signal) - 1)
k_arr = range(-len(ocitanje_signal)+1, len(referentni_genom_signal))

FFT_ref = np.fft.fft(padding_ref + referentni_genom_signal)
FFT_read = np.fft.fft(ocitanje_signal + padding_ocitanja)

correlation_freq = np.conjugate(FFT_read)*FFT_ref

cross_correlation = np.fft.ifft(correlation_freq)

correlation = np.abs(cross_correlation)

k = k_arr[cross_correlation.argmax()]

print(f"korelacija za očitavanje {i+1}:")
print(correlation)
print("Indeks najveće korelacije: {}".format(k))
print("-----")

```

```

korelacija za očitavanje 1:
[0.96628788 0.24537804 0.19311781 ... 1.28947659 1.7916225 0.59293484]
Indeks najveće korelacije: 2324486
-----

korelacija za očitavanje 2:
[1.66565335 0.37680388 2.8128802 ... 2.04208374 2.45570969 1.11019077]
Indeks najveće korelacije: 1877260
-----

korelacija za očitavanje 3:
[0.83991854 1.87479614 1.63159265 ... 3.14189966 1.76584122 0.7063178 ]
Indeks najveće korelacije: 557777
-----

korelacija za očitavanje 4:
[1.19272637 0.87937025 1.19654972 ... 0.79102689 0.39490803 0.79497562]
Indeks najveće korelacije: 1762444
-----

korelacija za očitavanje 5:
[1.2802567 0.24944829 2.69166605 ... 0.55792977 0.54486609 0.85331631]
Indeks najveće korelacije: 3583639
-----

korelacija za očitavanje 6:
[0.72018482 1.20032921 1.18881917 ... 2.70082707 1.66849989 0.79228674]

```

Indeks najveće korelacije: 4051518

korelacija za očitavanje 7:

[0.95922719 0.7894592 2.54336626 ... 1.42210027 0.0066199 0.63934382]

Indeks najveće korelacije: 2293706

korelacija za očitavanje 8:

[0.86276679 1.89785891 1.12991946 ... 1.53138375 0.8613733 0.95757401]

Indeks najveće korelacije: 1011323

korelacija za očitavanje 9:

[1.59628299 2.66051858 1.59656451 ... 1.42393059 2.07320064 0.82925831]

Indeks najveće korelacije: 628546

korelacija za očitavanje 10:

[1.14789333 1.9131893 0.76769862 ... 0.4511246 0.16230756 0.83104986]

Indeks najveće korelacije: 1497921

8. zadatak

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

```
[ ]: def hamming_distance(seq1, seq2):  
    assert len(seq1) == len(seq2), "Sekvence moraju biti iste duljine"  
    return sum(ch1 != ch2 for ch1, ch2 in zip(seq1, seq2))  
  
najveca_pozicija_korelacije = k  
  
# "Izrezujemo" dio referentne sekvence koji odgovara poziciji najveće korelacije  
pocetna_pozicija = najveca_pozicija_korelacije  
krajnja_pozicija = pocetna_pozicija + len(ocitanja_sekvenciranja[0].seq)  
referentni_segment = referentni_genom[0].seq[pocetna_pozicija:krajnja_pozicija]  
  
# Izračunavamo Hammingovu udaljenost između segmenta referentnog genoma i  
#   ↳ očitavanja  
udaljenost = hamming_distance(referentni_segment, ocitanja_sekvenciranja[0].seq)  
print("Hammingova udaljenost između segmenta referentnog genoma i prvog  
#   ↳ očitavanja je:", udaljenost)  
print(referentni_genom[0].seq[k:k+udaljenost])
```

```
print(ocitanja_sekvenciranja[0].seq[0:udaljenost])
```

Hammingova udaljenost između segmenta referentnog genoma i prvog očitavanja je: 9
GATCTGGTG
GATCTGGTG

```
[ ]: def hamming_distance(seq1, seq2):  
    assert len(seq1) == len(seq2), "Sekvence moraju biti iste duljine"  
    return sum(ch1 != ch2 for ch1, ch2 in zip(seq1, seq2))  
  
def najbolje_podudaranje(referenca, ocitanje):  
    min_udaljenost = float('inf')  
    najbolji_podniz = None  
    najveca_pozicija_korelacije = 0  
  
    for i in range(len(referenca) - len(ocitanje) + 1):  
        podniz = referenca[i:i+len(ocitanje)]  
        trenutna_udaljenost = hamming_distance(podniz, ocitanje)  
  
        if trenutna_udaljenost < min_udaljenost:  
            min_udaljenost = trenutna_udaljenost  
            najbolji_podniz = podniz  
            najveca_pozicija_korelacije = i  
  
    return najbolji_podniz, najveca_pozicija_korelacije  
  
# Pretpostavljamo da su referentni_genom i ocitanja_sekvenciranja već  
#   ↳ definirani negdje u kodu  
referentni_segment, pozicija = najbolje_podudaranje(referentni_genom[0].seq,  
#   ↳ ocitanja_sekvenciranja[0].seq)  
  
# Izračunavamo Hammingovu udaljenost između segmenta referentnog genoma i  
#   ↳ očitavanja  
udaljenost = hamming_distance(referentni_segment, ocitanja_sekvenciranja[0].seq)  
print("Hammingova udaljenost između segmenta referentnog genoma i prvog  
#   ↳ očitavanja je:", udaljenost)  
print("Najbolje podudaranje pronađeno na poziciji:", pozicija)
```

Hammingova udaljenost između segmenta referentnog genoma i prvog očitavanja je: 9
Najbolje podudaranje pronađeno na poziciji: 2324486

9. zadatak

U datoteci "ecoli_ILL_smallaln.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).

```
[ ]: import csv
```

```
tsv_file = open("./ecoli_ILL_small_aln.sam")
tsv_rows = csv.reader(tsv_file, delimiter="\t")
```

```
[ ]: header1, header2 = next(tsv_rows), next(tsv_rows)
first_record = next(tsv_rows)
print(first_record)
```

```
['SRR2052522.671', '0', 'NC_000913.3', '2324487', '60', '112M9S', '*', '0', '0',
'GATCTGGTGACCGGGTCGCGCAAAGTGATCATCGCCATGGAACATTGCGCCAAAGATGGTTCAGCAAAAATTTTGGGCC
TCTGTATCATGCCACTCACTGCGCAATATCCGGATCAAATGC', 'BGGG?D:GGDD=BFFBFFDEGDGDDGGDDDEBB
ED8DDF2??==GEDGDG<<>BGG>G>=A,A=8A<?1=???A#####
#####', 'NM:i:6', 'ms:i:188', 'AS:i:188', 'nn:i:0', 'tp:A:P', 'cm:i:13',
's1:i:78', 's2:i:0', 'dv:f:0.0320']
```

```
[ ]: print(f"Poravnanje definirano u 4. stupcu .sam datoteke : {first_record[3]}")
print(f"Moje poravnanje je {k}")
print(f"Razlika poravnanja jest, kao što je i navedeno : {int(first_record[3]) - k}")
```

Poravnanje definirano u 4. stupcu .sam datoteke : 2324487

Moje poravnanje je 2324486

Razlika poravnanja jest, kao što je i navedeno : 1

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.

```
[ ]: import Bio.pairwise2

korelacije = []

for i in range(100):
    ocitanje_signal = [enumerate_nucleotide(nucleotide) for nucleotide in
    ↪ ocitanja_sekvenciranja[i].seq]
    avg_ocitanja = np.average(ocitanje_signal)
    std_ocitanja = np.std(ocitanje_signal)
    ocitanje_signal = [(x - avg_ocitanja) / std_ocitanja for x in
    ↪ ocitanje_signal]

    padding_ocitanja = [0] * (len(referentni_genom_signal) - 1)
    padding_ref = [0] * (len(ocitanje_signal) - 1)
    k_arr = range(-len(ocitanje_signal)+1, len(referentni_genom_signal))

    FFT_ref = np.fft.fft(padding_ref + referentni_genom_signal)
    FFT_read = np.fft.fft(ocitanje_signal + padding_ocitanja)

    correlation_freq = np.conjugate(FFT_read)*FFT_ref

    cross_correlation = np.fft.ifft(correlation_freq)

    correlation = np.abs(cross_correlation)

    k = k_arr[cross_correlation.argmax()]

    korelacije.append(k)
```

```
[ ]: print(korelacije)
```

```
[2324486, 1877260, 557777, 1762444, 3583639, 4051518, 2293706, 1011323, 628546,
1497921, 3583853, 1212332, 1877622, 1890086, 767890, 2607362, 2279284, 2456798,
4434676, 3968956, 3382489, 4042795, 2098106, 832393, 1242185, 2171246, 147894,
1986927, 3802944, 3138889, 2487302, 2753343, 1441972, 4581225, 4403147, 1421664,
4092109, 772980, 3555436, 3307817, 2925637, 1190258, 4343721, 663438, 3251414,
3231183, 3426200, 4501246, 1287467, 1367056, 2909507, 4253639, 3508865, 61335,
344553, 401310, 2341959, 1630907, 3394077, 4632496, 4497880, 4224279, 3473354,
3547159, 1584282, 3443642, 2498150, 2349930, 4325020, 2359505, 3336147, 2545174,
220209, 227958, 2606563, 3973225, 3658914, 4356752, 1700776, 530761, 3201007,
3268804, 2718162, 1508527, 922585, 1094403, 3841115, 349911, 2101000, 4429448,
2705462, 4377909, 3667198, 3806074, 521347, 1219321, 3898429, 2487382, 3973268,
3733471]
```

```
[ ]: import pandas as pd

[ ]: with open("./ecoli_ILL_small_aln.sam") as tsv_file:
    tsv_rows = csv.reader(tsv_file, delimiter="\t")
    alignment_positions = [int(row[3]) for row in tsv_rows if not row[0].
        ↪startswith('@')]

count_within_5 = 0

for j in range(len(alignment_positions)):
    for i in range(len(korelacije)):
        difference = abs(alignment_positions[j] - korelacije[i])
        if difference <= 5:
            count_within_5 += 1

print(f"Broj očitavanja s razlikom manjom od 5: {count_within_5}")
```

Broj očitavanja s razlikom manjom od 5: 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.