

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 zvati *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.

```
In [ ]: %pip install biopython
```

Collecting biopython

Downloading biopython-1.81-cp311-cp311-win_amd64.whl (2.7 MB)

----- 2.7/2.7 MB 4.8 MB/s eta 0:00:00

Requirement already satisfied: numpy in c:\users\karlo\appdata\local\programs\python\python311\lib\site-packages (from biopython) (1.23.5)

Installing collected packages: biopython

Successfully installed biopython-1.81

Note: you may need to restart the kernel to use updated packages.

[notice] A new release of pip available: 22.3.1 -> 23.3.1

[notice] To update, run: python.exe -m pip install --upgrade pip

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.

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

```
import Bio
```

```
from Bio import SeqIO
```

```
import numpy as np
```

```
# Ispis verzije numpy
```

```
print("Verzija numpyja: ", np.__version__)
```

```
# Ispis verzije BioPython
```

```
print("Verzija biopythona: ", Bio.__version__)
```

```
Verzija numpyja: 1.23.5
```

```
Verzija biopythona: 1.81
```

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.

```
In [ ]: referencia = list(SeqIO.parse("escherichia_coli_reference.fasta", "fasta"))
print("Broj elemenata u referenci je: ", len(referenca))
ocitanja = list(SeqIO.parse("ecoli_ILL_small.fastq", "fastq"))
print("Broj ocitanja u listi je: ", len(ocitanja))

#S obzirom da je referencia samo jedna, mogu uzeti samo prvi element liste
referenca = referencia[0]
```

```
Broj elemenata u referenci je: 1
Broj ocitanja u listi 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 očitanje te identifikator i prvih 200 znakova za referentni genom E.coli.

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

```
In [ ]: print("Identifikator prvog očitanja E.coli: ", ocitanja[0].id, "sekvenca prvog očitanja: ", ocitanja[0].seq)
        print("Identifikator referentnog genoma E.coli: ", referenca.id, "prvih 200 znakova: ", referenca.seq[:200])

Identifikator prvog očitanja E.coli: SRR2052522.671 sekvenca prvog očitanja: GATCTGGTGACCGGGTCGCGCAAAGTGATCATCGCCATGGAACATT
GCGCCAAAGATGGTTCAGCAAAAATTTGGGCCTCTGTATCATGCCACTCACTGCGCAATATCCGGATCAAATGC
Identifikator referentnog genoma E.coli: NC_000913.3 prvih 200 znakova: AGCTTTTCATTCTGACTGCAACGGGCAATATGTCTCTGTGTGGATTAAAAA
AAGAGTGTCTGATAGCAGCTTCTGAAGTGGTTACCTGCCGTGAGTAAATTTAAATTTATTGACTTAGGTCACTAAATACTTTAACCAATATAGGCATAGCGCACAGACAGATAAAAAATTACAGA
GTACACAACATCCATGAAACGCAT
```

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.

```
In [ ]: def nukleotid_u_broj(nukleotid):
        vrijednosti={
            'A': 3,
            'G': 2,
            'C': -2,
            'T': -3
        }

        return vrijednosti.get(nukleotid.upper(), 0)
```

5. Zadatak

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

```
In [ ]: prvo_ocitanje = ocitanja[0]

ocit_prvo = [nukleotid_u_broj(nukleotid) for nukleotid in prvo_ocitanje.seq]
ref = [nukleotid_u_broj(nukleotid) for nukleotid in referenca.seq]

#print(prvo_ocitanje.seq, "\n", prvo_ocitanje_signal )
#C = np.correlate(prvo_ocitanje_signal, referenca_signal, 'full')
#print(C)

len1 = len(ref)
len2 = len(ocit_prvo)
k_arr = range(-len2 + 1, len1)
```

```

vrijednosti = [3, 2, -2, -3] # Vrijednosti za A, G, C, T da ne izvlacim iz funkcije od iznad

avg = np.average(vrijednosti)
std = np.std(vrijednosti)

ref = [(x - avg) / std for x in ref]
ocit_prvo = [(x - avg) / std for x in ocit_prvo]

padding1 = [0] * (len2 - 1)
padding2 = [0] * (len1 - 1)

X1 = np.fft.fft(padding1 + ref) # Referenca
X2 = np.fft.fft(ocit_prvo + padding2) # Očitanje

# Korelacija
Cor = np.conjugate(X2) * X1
cor = np.fft.ifft(Cor)
cor_real = np.real(cor)

# Maksimalna korelacija
k = k_arr[cor_real.argmax()]

# Ispis
print("Korelacija pomoću FFT (zanemarujući imaginarne vrijednosti):")
print(cor_real)
print("Maksimalna korelacija na poziciji k={}".format(k))

```

```

Korelacija pomoću FFT (zanemarujući imaginarne vrijednosti):
[-0.92307692  0.30769231 -0.15384615 ... -1.38461538 -1.84615385
 -0.61538462]
Maksimalna korelacija na poziciji 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.

```
In [ ]: print("Duljina reference: ", len(ref))

duljine_ocitanja = [len(ocitanje.seq) for ocitanje in ocitanja]
srednja_vrijednost = np.mean(duljine_ocitanja)
std_devijacija = np.std(duljine_ocitanja)

srednja_vrijednost, std_devijacija
```

Duljina reference: 4641652

```
Out[ ]: (121.0, 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.

```
In [ ]: for i in range(10): # Izračunajte korelaciju za prvih 10 očitavanja
    ocitanje = ocitanja[i]
    ocit = [nukleotid_u_broj(nukleotid) for nukleotid in ocitanje.seq]
    ocit = [(x - avg) / std for x in ocit]
    X2 = np.fft.fft(ocit + [0] * (len1 - 1)) # FFT za očitavanje s paddingom

    Cor = np.conjugate(X2) * X1
    cor = np.fft.ifft(Cor)
    cor_real = np.real(cor)
    k = k_arr[cor_real.argmax()]
    print(f"Korelacija za očitanje {i}:")
    print(cor_real)
    print(f"Maksimalna korelacija na poziciji k={k}")
```

Korelacija za očitavanje 0:

[-0.92307692 0.30769231 -0.15384615 ... -1.38461538 -1.84615385
-0.61538462]

Maksimalna korelacija na poziciji k=2324486

Korelacija za očitavanje 1:

[1.38461538e+00 -6.66134596e-15 -2.92307692e+00 ... -1.38461538e+00
-2.00000000e+00 -9.23076923e-01]

Maksimalna korelacija na poziciji k=1877260

Korelacija za očitavanje 2:

[-0.92307692 -2. -1.69230769 ... -2.76923077 -1.53846154
-0.61538462]

Maksimalna korelacija na poziciji k=557777

Korelacija za očitavanje 3:

[1.38461538 -0.46153846 -0.92307692 ... -1.38461538 -0.76923077
-0.92307692]

Maksimalna korelacija na poziciji k=1144877

Korelacija za očitavanje 4:

[1.38461538e+00 -8.43769769e-15 -2.46153846e+00 ... 1.53846154e-01
-7.69230769e-01 -9.23076923e-01]

Maksimalna korelacija na poziciji k=3583639

Korelacija za očitavanje 5:

[0.92307692 1.53846154 1.38461538 ... -3.23076923 -2.
-0.92307692]

Maksimalna korelacija na poziciji k=4051518

Korelacija za očitavanje 6:

[-9.23076923e-01 7.69230769e-01 2.46153846e+00 ... -1.38461538e+00
2.22046709e-15 6.15384615e-01]

Maksimalna korelacija na poziciji k=2293706

Korelacija za očitavanje 7:

[0.92307692 2. -1.07692308 ... 1.38461538 0.76923077
0.92307692]

Maksimalna korelacija na poziciji k=1011323

Korelacija za očitavanje 8:

[-1.38461538 -2.30769231 -1.38461538 ... -1.84615385 -2.30769231
-0.92307692]

Maksimalna korelacija na poziciji k=628546

Korelacija za očitavanje 9:

[-0.92307692 -1.53846154 0.92307692 ... -0.92307692 -0.46153846
-0.92307692]

Maksimalna korelacija na poziciji k=1497921

8. zadatak

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

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

```
def Hammingova_udaljenost(niz1, niz2):  
    if len(niz1) != len(niz2):  
        raise ValueError("Nizovi moraju biti iste duljine")  
    return sum(c1 != c2 for c1, c2 in zip(niz1, niz2))  
  
k = 2324486  
  
print("Hammingova razlika je ", Hammingova_udaljenost(prvo_ocitanje.seq, referencia.seq[k : k + len(prvo_ocitanje.seq)]))
```

Hammingova razlika je 9

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 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).

```
In [ ]: import csv
tsv_file = open("ecoli_ILL_small_aln.sam")
tsv_rows = csv.reader(tsv_file, delimiter="\t")

for row in tsv_rows:
    if row[0].startswith('@'):
        continue
    if row[0] == prvo_ocitanje.id:
        poravnanje_pozicija = int(row[3]) # Četvrti stupac sadrži poziciju
        break

# Usporedite pozicije
print("Pozicija dobivena korelacijom:", k) # k iz proslog zadatka
print("Pozicija iz SAM datoteke:", poravnanje_pozicija)
```

Pozicija dobivena korelacijom: 2324486

Pozicija iz SAM datoteke: 2324487

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.

```
In [ ]: from Bio import pairwise2

korelacije = []

def kalkulator_korelacije():
    for i in range(100):
        ocitanje = ocitanja[i]
        ocit = [nukleotid_u_broj(nukleotid) for nukleotid in ocitanje.seq]
        ocit = [(x - avg) / std for x in ocit]
        X2 = np.fft.fft(ocit + [0] * (len1 - 1)) # FFT za očitavanje s paddingom
        Cor = np.conjugate(X2) * X1
        cor = np.fft.ifft(Cor)
        cor_real = np.real(cor)
        k = k_arr[cor_real.argmax()]
        korelacije.append(k)

kalkulator_korelacije()
```

```
In [ ]: sam_korelacije = []

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

def kalkulator_sam_korelacije():
    i = 0
    for row in tsv_rows:
        if i >= 100:
            break

        if row[0].startswith('@'):
            continue
        if row[0] == ocitanja[i].id:
            sam_korelacije.append(int(row[3]))
```

```
i += 1
```

```
kalkulator_sam_korelacije()
```

```
In [ ]: print("Prvih 10 korelacija iz SAM filea: " , sam_korelacije[:10])
        print("Prvih 10 izračunatih korelacija: " , korelacije[:10])

        print("Broj slicnih je ", sum(abs(element1 - element2) <= 5 for element1, element2 in zip(korelacije, sam_korelacije)))

Prvih 10 korelacija iz SAM filea: [2324487, 1877261, 557778, 3910489, 20240, 4051519, 2808648, 1011324, 628547, 1497922]
Prvih 10 izračunatih korelacija: [2324486, 1877260, 557777, 1144877, 3583639, 4051518, 2293706, 1011323, 628546, 1497921]
Broj slicnih je 47
```

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.