

✓ HARINANDAN H

AM.EN.U4AIE22120

Problem 1

Generate a Random DNA Sequence Description: Create a random DNA string with letters from the whole alphabet A, C, G, and T. First make a list of random letters and then join all those letters to a string. Also write another function to count the number of bases in the random sequence and measure the CPU time for large such DNA strings. (Hint : use import random, import time)

```
import random
import time

bases = ['A', 'T', 'G', 'C']

dna_seq_len=10000
dna_seq=''.join(random.choice(bases)for i in range(dna_seq_len))
print(f'DNA_sequence:\n{dna_seq}')
print(f"\nA: {dna_seq.count('A')}")
print(f"T: {dna_seq.count('T')}")
print(f"G: {dna_seq.count('G')}")
print(f"C: {dna_seq.count('C')}")

start_time=time.process_time()
end_time=time.process_time()
cpu_time=end_time-start_time

print(f'\nCPU time: {cpu_time:.6f}seconds')

DNA_sequence:
TCCTCCCGCTTTTGGGAATGGATTAGGAGGGGATGCGGCCAATAAGCGATCTTTGACCATCTCCTAGCACACACGTGAACGTACGCGGCATACAATACACCGCGACGGGCGCGCGGGCGGTCGCTCGCAGGTATA

A: 2442
T: 2492
G: 2507
C: 2559

CPU time: 0.000030seconds
```

Problem 2:

Compute the Hamming Distance Between Two Strings We say that position i in k -mers $p_1 \dots p_k$ and $q_1 \dots q_k$ is a mismatch if $p_i \neq q_i$. For example, CGAAT and CGGAC have two mismatches. The number of mismatches between strings p and q is called the Hamming distance between these strings and is denoted $\text{HammingDistance}(p, q)$.

```
def Hamming_Distance(string_1, string_2):
    d=0
    if len(string_1)==len(string_2):
        for i in range(len(string_1)):
            if string_1[i]!=string_2[i]:
                d+=1
    return d

string_1="ACGTACGTA"
string_2="ACGTACGTT"

hamming_distance=Hamming_Distance(string_1, string_2)
print(hamming_distance)
```

1

Problem 3: Find Patterns Forming Clumps in a String Given integers L and t , a string Pattern forms an (L, t) -clump inside a (larger) string Genome if there is an interval of Genome of length L in which Pattern appears at least t times. For example, TGCA forms a $(25,3)$ -clump in the following Genome: gatcagcataagggtcccgTGAATGCATGACAAGCCTGCAGttgtttac Clump Finding Problem Find patterns forming clumps in a string. Given: A string Genome, and integers k , L , and t . Return: All distinct k -mers forming (L, t) -clumps

Problem 4: Find a Position in a Genome Minimizing the Skew Define the skew of a DNA string Genome, denoted $\text{Skew}(\text{Genome})$, as the difference between the total number of occurrences of 'G' and 'C' in Genome. Let $\text{Prefix}_i(\text{Genome})$ denote the prefix (i.e., initial substring) of Genome of length i . For example, the values of $\text{Skew}(\text{Prefix}_i(\text{"CATGGGCATCGGCCATACGCCCATGGGCATCGGCCATACGCC"}))$ are: 0 -1 -1 -1 0 1 2 1 1 0 1 2 1 0 0 0 -1 0 -1 -2 Minimum Skew Problem Find a position in a genome minimizing the skew. Given: A DNA string Genome. Return: All integer(s) i minimizing $\text{Skew}(\text{Prefix}_i(\text{Text}))$ over all values of i (from 0 to $|\text{Genome}|$).

8

