# DNA Sequence Analyzer Using Python

Author: JYOTSNA.R

Institution: Rajalakshmi Engineering College

Year: 2025

## Abstract

This project presents a basic DNA Sequence Analyzer built using Python. It is designed for biotechnology students with beginner-level programming experience. The tool performs key DNA analyses including GC content calculation, codon frequency detection, reverse complement generation, and motif identification. This tool helps visualize and understand biological sequence data and can be expanded for real-world datasets and research.

## Objective

- To develop a beginner-friendly Python tool for analyzing DNA sequences.

- To understand and implement bioinformatics concepts such as GC content, motifs, codons, and reverse complements.

- To demonstrate the practical use of Python in biotechnology.

## Background

DNA sequences are fundamental to genetics and bioinformatics. Simple sequence analysis can reveal important insights such as gene structure, stability (GC content), translation start points (codons), regulatory regions (motifs), and the complementary strand information.

## Key Concepts

- GC Content: Proportion of guanine (G) and cytosine (C) bases in DNA; higher GC implies more stability.

- Codons: Groups of three bases representing an amino acid or control signal.

- Motifs: Short, biologically significant patterns like "TATA" (promoter regions).

- Reverse Complement: DNA strand opposite to the given strand in 5' to 3' direction.

## Tools Used

- Python 3.x

- Libraries: collections, matplotlib (for visualization)

## Sample DNA Sequence

sequence = "ATGCTATATAGCCGATATAAGG"

## Step 1: GC Content Calculation

```
def gc_content(seq):

    gc = seq.count('G') + seq.count('C')

    return (gc / len(seq)) * 100
```

## Step 2: Most Frequent Codon

```
from collections import Counter


def most_frequent_codon(seq):

    codons = [seq[i:i+3] for i in range(0, len(seq)-2, 3)]

    return Counter(codons).most_common(1)[0]
```

## Step 3: Reverse Complement

```
def reverse_complement(seq):

    complement = {'A': 'T', 'T': 'A', 'G': 'C', 'C': 'G'}

    rev_comp = ''.join([complement[base] for base in reversed(seq)])

    return rev_comp
```

## Step 4: Motif Search

```
def find_motif(seq, pattern):

    positions = []

    for i in range(len(seq) - len(pattern) + 1):
```

```
        if seq[i:i+len(pattern)] == pattern:

            positions.append(i)

    return positions
```

## Step 5: GC Content Visualization

```
import matplotlib.pyplot as plt


def gc_content_window(seq, window=5):

    gc_vals = [gc_content(seq[i:i+window]) for i in range(0, len(seq) - window + 1)]

    plt.plot(gc_vals)

    plt.title("GC Content Across DNA")

    plt.xlabel("Position")

    plt.ylabel("GC %")

    plt.show()
```

## Results and Output

- GC Content: ~45.45%

- Most frequent codon: ATA (3 times)

- Reverse complement: CCTTATATCGGCTATATAGCAT

- Motif "TATA" found at positions: 4, 6, 15

- GC content visualized as a graph

## Optional Enhancements

- Accept DNA input from users

- Support for FASTA files

- Codon to amino acid translation

- Graphical motif mapping

## Conclusion

This project demonstrates how basic bioinformatics tasks can be performed using Python. It provides a foundation for students in biotechnology to build analytical tools for sequence data. With additional features, this tool can evolve into a more comprehensive DNA analysis pipeline.

## Future Work

- Apply the tool to real gene sequences from NCBI databases

- Add machine learning for gene prediction

- Build a web interface using Flask or Streamlit

## References

- Bioinformatics Algorithms by Phillip Compeau and Pavel Pevzner

- NCBI GenBank

- Python documentation