

Workshop Report 1: Systems Analysis and Motive Detection

Introduction

This report presents a novel Java-based algorithm for identifying recurring patterns, or motifs, within artificially generated genetic sequences. The algorithm was developed as part of the first workshop of the Systems Analysis course. By incorporating entropy-based filtering, the algorithm effectively distinguishes meaningful motifs from random occurrences. Our experimental results demonstrate the algorithm's accuracy and efficiency in detecting motifs of varying lengths, providing valuable insights into the challenges and complexities of motif identification in genetic sequences.

Objective

The main objective of the program is to simulate a database of artificial genetic sequences and detect the motif (character pattern) that appears most frequently within a specific size. Additionally, an entropy-based filter is implemented to select more diverse and potentially chaotic sequences.

Systemic Analysis

The program is made up of three main modules:

1. **Generation of Artificial Sequences:** This module creates a set of n random genetic sequences of size m . Each sequence is composed of nucleotides represented by the letters A, C, G and T. The selection probabilities of each nucleotide are configurable parameters. The generated sequences are saved in a text file.

```
1 private static String generarSecuencia(int tamaño, double probabilidadA, double probabilidadC, double probabilidadG, double probabilidadT) {
2     // Genera una secuencia genética artificial de tamaño 's' con las probabilidades dadas
3     StringBuilder sb = new StringBuilder();
4     Random random = new Random();
5     for (int i = 0; i < tamaño; i++) {
6         double aleatorio = random.nextDouble();
7         if (aleatorio < probabilidadA) {
8             sb.append('A');
9         } else if (aleatorio < probabilidadA + probabilidadC) {
10            sb.append('C');
11        } else if (aleatorio < probabilidadA + probabilidadC + probabilidadG) {
12            sb.append('G');
13        } else {
14            sb.append('T');
15        }
16    }
17    return sb.toString();
18 }
```

```

1 // Crear secuencias genéticas artificiales
2 String[] secuencias = new String[númeroSecuencias];
3 for (int i = 0; i < númeroSecuencias; i++) {
4     secuencias[i] = generarSecuencia(tamañoSecuencia, probabilidadA, probabilidadC, probabilidadG, probabilidadT);
5 }

```

2. **Entropy Calculation and Filtering:** This module calculates the Shannon entropy for each sequence. Entropy is a measure of disorder or chaos in a system. In this case, it is used to filter sequences with low entropy, that is, with high repetition of the same nucleotide. Sequences with entropy greater than or equal to the average entropy are retained for subsequent analysis.

```

1 private static double calcularEntropía(String secuencia) {
2     // Calcular la entropía de una secuencia genética
3     int[] frecuencias = new int[4]; // A, C, G, T
4     for (char nucleótido : secuencia.toCharArray()) {
5         switch (nucleótido) {
6             case 'A' -> frecuencias[0]++;
7             case 'C' -> frecuencias[1]++;
8             case 'G' -> frecuencias[2]++;
9             case 'T' -> frecuencias[3]++;
10        }
11    }
12    double entropía = 0;
13    for (int frecuencia : frecuencias) {
14        if (frecuencia > 0) {
15            double probabilidad = (double) frecuencia / secuencia.length();
16            entropía -= probabilidad * Math.log(probabilidad) / Math.log(2);
17        }
18    }
19    return entropía;
20 }
21 }

```

```

1 // Calcular entropía promedio
2 double sumaEntropía = 0;
3 for (String secuencia : secuencias) {
4     sumaEntropía += calcularEntropía(secuencia);
5 }
6 double entropíaPromedio = sumaEntropía / secuencias.length;
7

```

3. **Motif Detection:** This module searches for the motif of size s that is most frequently repeated within the filtered sequences. All possible combinations of nucleotides of size s in each sequence are evaluated. If there are patterns with the same frequency, the one with the greatest number of consecutive repeat bases is selected.

```
1 // Leer archivo y buscar la secuencia de caracteres de tamaño s que más se repite
2 Map<String, Long> frecuenciasSecuencias = new HashMap<>();
3 long startTime = System.nanoTime(); // Comenzar el cronometraje
4 try (BufferedReader reader = new BufferedReader(new FileReader(archivoSalidaFiltradas))) {
5     String linea;
6     while ((linea = reader.readLine()) != null) {
7         for (int i = 0; i <= linea.length() - tamañoSecuenciaBuscar; i++) {
8             String secuencia = linea.substring(i, i + tamañoSecuenciaBuscar);
9             frecuenciasSecuencias.put(secuencia, frecuenciasSecuencias.getOrDefault(secuencia, 0L) + 1);
10        }
11    }
```

```
1 // Encontrar secuencia más frecuente
2 String secuenciaMasFrecuente = null;
3 long frecuenciaMaxima = 0;
4 for (Map.Entry<String, Long> entry : frecuenciasSecuencias.entrySet()) {
5     if (entry.getValue() > frecuenciaMaxima) {
6         frecuenciaMaxima = entry.getValue();
7         secuenciaMasFrecuente = entry.getKey();
8     }
9 }
10
```

Complexity

The generation of artificial sequences has a time complexity of $O(n \times m)$, where n is the number of sequences and m is the size of each sequence.

The entropy calculation for each sequence has a time complexity of $O(m)$, where m is the size of the sequence.

Entropy-based sequence filtering has a time complexity of $O(n)$, where n is the total number of sequences.

Motif detection has a worst-case time complexity of $O(n \times m \times s^3)$, where n is the number of filtered sequences, m is the size of the filtered sequences, and s is the size of the searched motif. This is because it iterates over all sequences, analyzes substrings

of size s , and in the worst case all possible combinations of nucleotides for the motif (s^3) must be evaluated.

Chaos Analysis

Entropy is used as an indirect measure of chaos in the system. Sequences with high entropy have greater diversity in their nucleotides, which can be considered a more chaotic system. The entropy-based filter allows selecting sequences that potentially contain less repetitive and possibly more interesting patterns for motif analysis.

Results

The program is executed taking as parameters the number of sequences (n), the size of the sequences (m), the probabilities of each nucleotide, the size of the motif to be searched (s) and the names of the files to save the generated sequences. and the filtered sequences.

The program output shows:

- The most frequent motif of size s found in the filtered sequences.
- The frequency of appearance of the reason.
- The average entropy of the generated sequences.
- The search time for the reason in seconds.
- The total execution time of the program in seconds.

Discussion of Results

The results of the program largely depend on the configured parameters. The selection of nucleotide probabilities and the size of the motif to search will influence the diversity of the sequences and the ease of detecting patterns. The average entropy can give an idea of the level of chaos in the generated sequences, and the execution time allows evaluating the efficiency of the algorithm used for motif detection.

Conclusion

The program implements a basic approach for motif detection in artificial genetic sequences. Random generation, entropy calculation and exhaustive search techniques are used to achieve the objective. Complexity analysis allows you to understand program performance based on the amount of data and the size of the desired motif. Incorporating entropy as a filter allows the analysis to be biased towards sequences with greater potential diversity.