

Instituto Federal Catarinense - Campus Videira

Aluno: Dyeizon Procopiuk da Silva

Turma: Ciência da Computação 2022/1

Disciplina: Algoritmos

Data: 02/07/2022

Metodologia utilizada para a solução do problema

Para a entrada da sequência de DNA, foi criado um arquivo `.txt` que deve conter na primeira linha a dimensão da matriz quadrada, e nas linhas seguintes, a matriz quadrada a ser lida. Cada elemento da matriz deve estar sendo separado por um caractere de espaço. Então, a função **getDNAFromFile** interpreta o arquivo passado como parâmetro, e aplica as alterações necessárias para que a matriz seja utilizada pelo algoritmo, retornando um vetor com os elementos concatenados, resultando em uma sequência de DNA crua.

Figura 1 - Modelo de entrada

```
1 20
2 C T A T T T C C T G C T T G C A G C G
3 G T G T T A T A A C A A C C C T C C C T
4 T T T T A G C C C T C G C T A T G T A G
5 G T A A A C A A A G T A G T G T T A C C
6 A T C A C C A G G C A G C A G C C C G A
7 A A T A C G C T A T G T G T A A A T T C
8 A T C T T G A T C A C C C A C T G C C T
9 A G C C A T A G A A A C G A T T G T C G
10 T A G A A T A C A A A A G C C G G T A C
11 G C T C A G T G A T T T A G G T T A G G
12 C G T A A A A C G A T C G A T A A T C T
13 A T A C G T T G T G T G T A C G T A A T
14 T A C C T C C A G T G A T T A C C C G A
15 G G T C T C T A A T A A T A A G A G A T
16 T C C A G C G A C A G T A T G C A T T G
17 G C A G C T T T C C C C C T A C A A G T
18 T G G T A T T T A A G A A G G C G A G G
19 A C C G G T A A A A T A C C T G T T G T
20 T G A A T T G G C A G T G T C C A G C T
21 G A C A G T A T C C C G A C T T C C C G
```

O próximo passo é a validação da sequência de DNA. A função **checkDNA** verifica se há alguma base nitrogenada inválida dentro da sequência de DNA, considerando que uma base nitrogenada inválida é qualquer caractere que não seja "A", "G", "C" ou "T".

Após a validação da sequência de DNA, é criada uma instância do registro **Especime** (Figura 2), utilizando alocação dinâmica de memória. Este registro armazena a sequência de DNA, e a espécie que esta sequência corresponde.

Figura 2 - Registro Especime

```
typedef struct Especime {  
    char dna;  
    char tipo[20];  
} Especime;
```

A função **isSimian** submete a sequência de DNA a uma série de testes nas linhas, colunas, diagonais principais e diagonais secundárias em busca de repetições de 4 bases nitrogenadas coincidentes. Por motivos de performance, é verificado primeiramente as linhas e colunas, tendo a matriz inteira como referência, e caso seja encontrado repetições, a função responsável por realizar os testes das diagonais não é chamada.

Caso não sejam encontradas repetições nas linhas e colunas, a função recursiva **subMatrix** é chamada. Esta função é responsável por procurar repetições na diagonal principal e secundária de uma submatriz 4x4 que percorrerá toda a matriz (ambas as verificações são feitas ao mesmo tempo na submatriz), e chamando a si própria com valores diferentes até encontrar repetições, ou até chegar ao final da matriz.

O retorno da função **isSimian** determinará o tipo da espécie analisada pela sequência de DNA, sendo “humano” caso retorne 0, e “símio” caso retorne 1. Caso seja um símio, o algoritmo informa ao usuário onde foi encontrada repetição e qual caractere foi encontrado repetidamente.

O algoritmo pode ser facilmente modificado para considerar repetições mais extensas de bases nitrogenadas, assim como maiores submatrizes para as verificações das diagonais, caso necessário.

Resultados obtidos

Realizei uma sequência de testes com uma matriz quadrada de dimensões 20x20, em que inseri repetições de bases nitrogenadas em posições aleatórias da matriz, para verificar a precisão do algoritmo.

A primeira verificação foi nas linhas (Figura 3), o algoritmo encontrou uma repetição de 4 bases nitrogenadas “A” na linha 17 em 0.000238 segundos. O algoritmo finalizou a execução no momento em que a repetição foi encontrada, e **não** verificou as colunas e nem as diagonais, pois o fato de haver ou não repetições de bases nitrogenadas em outros sentidos não alteraria o resultado.

A próxima verificação é a das colunas (Figura 4), o algoritmo encontrou a repetição da base nitrogenada “G” em 0.000267 segundos.

Ao não encontrar repetições de bases nitrogenadas nas linhas e colunas, a função **subMatrix** é chamada para tentar localizar repetições nas diagonais. No próximo exemplo, foi encontrada uma repetição da base nitrogenada “C” na diagonal principal em 0.000330 segundos (Figura 5).

A última verificação da **subMatrix** é na diagonal secundária, para este teste em questão, coloquei a repetição na última posição possível da submatriz. A repetição da base nitrogenada “T” foi encontrada na diagonal secundária 0.000356 segundos (Figura 6).

Figura 3 - Linhas

```

A T A T T A T C C T G C T T G C A G C G
G T G A T A T A A C A A C C C T C C C T
T A T T A G C C C T C G C T A T G T A G
G T A T A C C A A G T A G T G T T A C C
A T C A A T A C G C A G C A G C C C G A
T A T A C G C T C C G T G T A A A T T C
A T C T A G A T C C A C C A C T G C C T
A G C C A A A G A A A C G A T T G T C G
T A G A A T T C A A G A T C C G G T A C
G C T C G G T T A T T T G G G T T A G G
C G T G A A A C G A T C T A T A A T C T
A T A C G T T G T G T G G A G G T A A T
T A C C T C C A G T G A T T G C C C G A
G G T C T C T A A T A A T A A G A G A T
T C C A G C G A C A G T A T G C A T T G
G C A G C G T T C C G C T T A C G A G T
T G G T A T G T A A G T A G G C A A G C
A C C G G G T A A A A G C C T G T G A T
T G A A T T G G C A G T G T C C A C G T
G A C A G T A T C C C G A C T T C C C G
Sequência de A na linha 17
É um símio
Tempo de execução: 0.000238s

```

Figura 4 - Colunas

```

A T A T T A T C C T G C T T G C A G C G
G T G A T A T A A C A A C C C T C C C T
T A T T A G C C C T C G C T A T G T A G
G T A T A C C A A G T A G T G T T A C C
A T C A A T A C G C A G C A G C C C G A
T A T A C G C T C C G T G T A A A T T C
A T C T A G A T C C A C C A C T G C C T
A G C C A A A G A A A C G A T T G T C G
T A G A A T T C A A G A T C C G G T A C
G C T C G G T T A T T T G G G T T A G G
C G T G A A A C G A T C T A T A A T C T
A T A C G T T G T G T G G A G G T A A T
T A C C T C C A G T G A T T G C C C G A
G G T C T C T A A T A A T A A G A G A T
T C C A G C G A C A G T A T G C A T T G
G C A G C G T T C C G C T T A C G A G T
T G G T A T G T A A G T A G G C A A G C
A C C G G G T A T A A G C C T G T G A T
T G A A T T G G C A G T G T C C A C G T
G A C A G T A T C C C G A C T T C C C G
Sequência de G na coluna 14
É um símio
Tempo de execução: 0.000267s

```

Figura 5 - Diagonais principais

```

A T A T T A T C C T G C T T G C A G C G
G T G A T A T A A C A A C C C T C C C T
T A T T A G C C C T C G G T A T G T A G
G T A T A C C A A G T A C T G T T A C C
A T C A A T A T G C A G C A G C C C G A
T A T A C G C T C T G T G T C A A T T C
A T C T A G A T C C A C C A C C G C C T
A G C C A A A G A A A C G A T T G T C G
T A G A A T T C A A G A T C C G G T A C
G C T C G G T T A T T T G G G T T A G G
C G T G A A A C G A T C T A T A A T C T
A T A C G T T G T G T G G A G G T A A T
T A C C T C C A G T G A T T G C C C G A
G G T C T C T A A T A A T A A G A G A T
T C C A G C G A C A G T A T G C A T T G
G C A G C G T T C C G C T T A C G A G A
T G G T A T G T A A G T A G G C A A G T
A C C G G G T A T A A G C C T G T G T T
T G A A T T G G C A G T G T C C A T G T
G A C A G T A T C C C G A C T T A C C G
Sequência de C na diagonal principal,
posição final: [6, 15]
É um símio
Tempo de execução: 0.000330s

```

Figura 6 - Diagonais secundárias

```

A T A T T A T C C T G C T T G C A G C G
G T G A T A T A A C A A C C C T C C C T
T A T T A G C C C T C G G T A T G T A G
G T A T A C C A A G T A C T G T T A C C
A T C A A T A T G C A G C A G C C C G A
T A T A C G C T C T G T G T C A A T T C
A T C T A G A T C C A C C A C C G C C T
A G C C A A A G A A A C G A T T G T C G
T A G A A T T C A A G A T C C G G T A C
G C T C G G T T A T T T G G G T T A G G
C G T G A A A C G A T C T A T A A T C T
A T A C G T T G T G T G G A G G T A A T
T A C C T C C A G T G A T T G C C C G A
G G T C T C T A A T A A T A A G A G A T
T C C A G C G A C A G T A T G C A T T G
G C A G C G T T C C G C T T A C G A G A
T G G T A T G T A A G T A G G C A A G T
A C C G G G T A T A A G C C T G T G T T
T G A A T T G G C A G T G T C C A T G T
G A C A G T A T C C C G A C T T A C C G
Sequência de T na diagonal secundária,
posição final: [19, 16]
É um símio
Tempo de execução: 0.000356s

```

Caso sejam feitas todas as verificações e não tenham sido encontradas repetições, o espécime é considerado “humano”, como verificado na Figura 7, em que a verificação levou 0.000390 segundos para passar por todas as etapas de verificação, e chegar à conclusão de que a sequência de DNA corresponde a de um ser humano.

Figura 7 - Sequência de DNA humano

```
A T A T T A T C C T G C T T G C A G C G
G T G A T A T A A C A A C C C T C C C T
T A T T A G C C C T C G G T A T G T A G
G T A T A C C A A G T A C T G T T A C C
A T C A A T A T G C A G C A G C C C G A
T A T A C G C T C T G T G T C A A T T C
A T C T A G A T C C A C C A C C G C C T
A G C C A A A G A A A C G A T T G T C G
T A G A A T T C A A G A T C C G G T A C
G C T C G G T T A T T T G G G T T A G G
C G T G A A A C G A T C T A T A A T C T
A T A C G T T G T G T G G A G G T A A T
T A C C T C C A G T G A T T G C C C G A
G G T C T C T A A T A A T A A G A G A T
T C C A G C G A C A G T A T G C A T T G
G C A G C G T T C C G C T T A C G A G A
T G G T A T G T A A G T A G G C A A G T
A C C G G G T A T A A G C C T G T G T T
T G A A T T G G C A G T G T C C A T G T
G A C A G T A T C C C G A C T T G C C G
É um humano
Tempo de execução: 0.000390s
```