



Facultad de Matemática y Computación
Universidad de la Habana

Proyecto de Programación I Moog!e!

Ernesto Abreu Peraza
Grupo: C-121
Curso: 2023

Descripción

Moogle! es una aplicación *"totalmente original"* cuyo propósito es buscar inteligentemente un texto en un conjunto de documentos.

Es una aplicación web, desarrollada con tecnología .NET Core 6.0, específicamente usando Blazor como framework web para la interfaz gráfica, y en el lenguaje C#. La aplicación está dividida en dos componentes fundamentales:

- [MoogleServer](#) es un servidor web que renderiza la interfaz gráfica y sirve los resultados.
- [MoogleEngine](#) es una biblioteca de clases donde está... ehem... casi implementada la lógica del algoritmo de búsqueda.

Correr y usar el proyecto

Para correr el proyecto debes usar el comando `dotnet watch run -project MoogleServer` en Windows y `make dev` en Linux. En la carpeta [Content](#) deberán aparecer los documentos (en formato *.txt) en los que el usuario va a realizar la búsqueda. En la casilla donde aparece *Introduzca la búsqueda* el usuario va a escribir que desea buscar y basta con apretar el botón *Buscar* para que Moogle! haga su trabajo.

Arquitectura del proyecto

Características de clases

- [Moogle](#): procesar consulta.
- [TF_IDF](#): calcular tf-idf para los documentos y para la query, calcular coseno del angulo entre dos vectores.
- [DocumentReader](#): leer, normalizar y obtener de los documentos, título, texto y palabras.
- [StringUtil](#): procesar, calcular, modificar y obtener informacion a partir de texto util para alguna de las funcionalidades de la aplicación.
- [Matrix](#): Definir y multiplicar matrices
- [Vector](#): Definir, calcular modulo y multiplicar vectores

Precálculo

Term frequency-Inverse document frequency (TF_IDF) es una medida numérica que expresa cuan relevante es una palabra para un documento en una colección de documentos.

$$TF_IDF[t, d] = TF[t, d] * IDF[t]$$
$$TF[t, d] = \frac{f[t]}{maxFrequency}$$
$$IDF[t] = \log_{10} \frac{numberOfDocuments}{Df[t]}$$

Siendo t una palabra y d un documento. $f[t]$ la frecuencia de la palabra t en el documento d y $maxFrequency$ es el máximo de los $f[t]$. $numberOfDocuments$ el total de documentos y $Df[t]$ es la cantidad de documentos en los que aparece la palabra t .

- Cuando el programa empieza a correr se ejecuta el método `TF_IDF.Compute()` del archivo `TF_IDF.cs`. Este método lee el texto de cada documento que aparece en la carpeta `Content` y de él extrae todas las palabras. Se calcula el TF_IDF para cada palabra en cada documento

```
1 public static void Compute()
2 {
3     /* Precalcula el TF_IDF */
4
5     /* Arreglo con los nombres de los documentos */
6     documentsName = DocumentReader.DocumentsNameList("../\\Content");
7
8     int numberOfDocuments = documentsName.Length;
9
10    /* Diccionario [nombre de documento => indice] */
11    for (int i = 0; i < numberOfDocuments; i++)
12    {
13        Document[documentsName[i]] = i;
14    }
15
16    /* Arreglo con las palabras de los documentos */
17    words = DocumentReader.WordsList(documentsName);
18
19    int numberOfWords = words.Length;
20
21    /* Diccionario [palabra => indice] */
22    for (int i = 0; i < numberOfWords; i++)
23    {
24        Word[words[i]] = i;
25    }
26
27    Matrix TF = ComputeTF();
28    tf = TF;
29
30    Vector IDF = ComputeIDF();
31    idf = IDF;
32
33    Matrix TF_IDF = new Matrix(numberOfWords, numberOfDocuments);
34
35    /* Multiplicar TF por IDF */
36    for (int j = 0; j < numberOfDocuments; j++)
37        for (int i = 0; i < numberOfWords; i++)
38            TF_IDF[i, j] = TF[i, j] * IDF[i];
39
40    tf_idf = TF_IDF;
41 }
```

- Los métodos `ComputeTF()` y `ComputeIDF()` calculan el TF y el IDF respectivamente.

Consulta

- Luego que la aplicación inicie, cuando se realice una consulta (búsqueda), el proyecto llama al método `Moogle.Query()` del archivo `Moogle.cs` donde se calcula el TF_IDF para la consulta a través del método `TF_IDF.ComputeQueryTF_IDF()` y se compara con el TF_IDF de cada documento dándole una puntuación a cada uno y devolviendo una lista con los nombres de los documentos que más relevancia tienen. La puntuación sería el coseno del ángulo entre el vector formado con el TF_IDF de la consulta y el TF_IDF del documento. La similitud de los documentos depende es mayor mientras mas se acerca a 1 el coseno del angulo entre ellos, visto de otra forma mientras mas pequeño es el angulo entre ambos. Para esto se utilizan dos fórmulas de dot product. Para dos vectores a y b: $dotProduct = a_1 * b_1 + a_2 * b_2 + ... + a_n * b_n$, $dotProduct = \cos(a, b) * |a| * |b|$.

```
1 public static (float, int)[] VectorialModel(Vector queryTF_IDF)
2 {
3     /* Devuelve un arreglo que contiene para cada documento el coseno
4     del angulo
5     entre el vector de la query y el del documento */
6     (float, int)[] vectorialModel = new (float, int)[tf_idf.columns];
7     for (int j = 0; j < tf_idf.columns; j++)
8     {
9         Vector v = new Vector(tf_idf.rows);
10        for (int i = 0; i < tf_idf.rows; i++)
11        {
12            v[i] = tf_idf[i, j];
13        }
14        float score = 0;
15        if (queryTF_IDF.Module() * v.Module() != 0)
16            score = Vector.Dot_Product(queryTF_IDF, v) / (queryTF_IDF.Module
17            () *
18            v.Module());
19        vectorialModel[j] = (score, j);
20    }
21    return vectorialModel;
22 }
23 ...
24 /* Definicion de producto escalar entre dos vectores */
25 static public float Dot_Product(Vector a, Vector b)
26 {
27     if (a.Dimensions != b.Dimensions)
28         /* Exception */
29         return 0;
30
31     float dot_product = 0;
32     for (int i = 0; i < a.Dimensions; ++i)
33         dot_product += a[i] * b[i];
34
35     return dot_product;
36 }
37 ...
38 ...
39 ...
40 /* Definicion de modulo de un vector */
```

```
41 public float Module()
42 {
43     float module = 0;
44     for (int i = 0; i < this.Dimensions; i++)
45     {
46         module += this[i] * this[i];
47     }
48     return (float)Math.Sqrt(module);
49 }
```

- También se muestra un fragmento por cada documento donde se puede apreciar en este una relación del documento con la consulta.

Funcionalidades extras

Otra funcionalidad del Moogle! es que dará una sugerencia de búsqueda en caso de que el usuario quizás cometió un error al escribir la consulta. Para esto usamos un algoritmo de Edit Distance conocido como Levenshtein distance y a este costo le dividimos por el Longest Common Prefix. El propósito de este último proviene de la idea de que es más probable equivocarse en las últimas letras que en las primeras. Así distra está más cerca de dijkstra que de citra.

```
1 public static float Distance(string a, string b)
2 {
3     /* Calcula la distancia entre dos cadenas de caracteres */
4     return EditDistance(a, b) / LongestCommonPrefix(a, b);
5 }
6
7 public static float EditDistance(string a, string b)
8 {
9     /* Calcula el EditDistance para dos cadenas de caracteres */
10    int[,] dp = new int[a.Length + 1, b.Length + 1];
11    for (int i = 0; i <= a.Length; i++)
12        for (int j = 0; j <= b.Length; j++)
13        {
14            if (i == 0 || j == 0)
15            {
16                dp[i, j] = Math.Max(i, j);
17            }
18            else
19            {
20                dp[i, j] = Int32.MaxValue;
21                if (a[i - 1] == b[j - 1])
22                    dp[i, j] = dp[i - 1, j - 1];
23                else
24                {
25                    dp[i, j] = Math.Min(dp[i, j], dp[i - 1, j] + 1);
26                    dp[i, j] = Math.Min(dp[i, j], dp[i, j - 1] + 1);
27                    dp[i, j] = Math.Min(dp[i, j], dp[i - 1, j - 1] + 1);
28                }
29            }
30        }
31    return (float)dp[a.Length, b.Length];
32 }
33
```

```
34 public static float LongestCommonPrefix(string a, string b)
35 {
36     /* Calcula el LongestCommonPrefix para dos cadenas de caracteres */
37     for (int i = 0; i < Math.Min(a.Length, b.Length); i++)
38         if (a[i] != b[i])
39             return (float)(i + 1);
40
41     return (float)Math.Min(a.Length, b.Length) + 1;
42 }
```