

Moogle! es un programa que hace búsquedas en tiempo real sobre una base de datos de documentos.

En principio necesita una base de datos de archivos en formato .txt ubicados en la carpeta Content.

Cuenta además del archivo principal Moogle.cs con 5 archivos .cs necesarios (bueno 4 necesarios para funcionar) que son):

- AuxItems, una clase que guarda el texto y título de un documento.
- Methods, una clase que tiene todos los métodos que son invocados en el programa.
- SearchItem, una clase que guarda el título, snippet y score de cada documento al hacer la query
- SearchResult, que como su nombre indica devuelve el resultado de la búsqueda.
- Matrix, una clase que ayuda en el trabajo con operaciones con matrices, pero que no fue necesario usar en el funcionamiento del proyecto.

Lo primero que hace el proyecto es hacer un arreglo de AuxItems donde en cada posición hay un documento de la carpeta Content.

Luego duplica ese arreglo pero con los documentos normalizados, esto es, los mismos textos de los documentos, pero sin signos de puntuación, mayúsculas ni saltos de línea.

Ahora creo un diccionario y voy recorriendo los textos de mis documentos. Cada vez que encuentro una palabra "nueva" para mi diccionario la anado con un número natural que significa que número de palabra es.

Paralelo a esto voy creando una matriz de (palabras;documentos).

Al terminar esta parte del código me queda un diccionario en el que todas las palabras presentes se les asigna un número y una matriz de tamaño [cant de palabras; cant de documentos]

Aquí empieza la parte algebraica, pues por cada documento voy calculando la importancia que tiene cada una de las palabras que lo forman, a través de un método llamado TF\*IDF. En este punto ya tengo el IDF de cada palabra de mi universo ya calculada (ya que este no varía) y solo llamo al método privado TF que calcula el TF de cada una de las palabras de un documento y al multiplicarlas por su respectivo IDF me queda la importancia de esas palabras para ese documento. Luego le asigno ese valor a su posición correspondiente en la matriz.

Al acabar me queda la matriz con el TF\*IDF de cada palabra de mi universo de palabras en cada documento.

En este punto del programa es que valoro la query. Lo primero que hago es normalizarla.

Luego le calculo el TF\*IDF para ver que tan relevante es para mi universo. En caso de que no tenga palabras en común con mis documentos su relevancia va a ser obviamente 0, y a partir de ahí se incrementará en dependencia.

Una vez que tengo el TF\*IDF de la query puedo aplicar una fórmula llamada "Similitud Coseno" que me dice el ángulo de desviación de un vector respecto a otro. Si ese coseno da 0 entonces ambos vectores apuntan al mismo lugar.

Aplicando esto al proyecto. Si trato los documentos como un conjunto de vectores y la query como otro la similitud de coseno me dirá que tan parecidos son. Dicho y hecho: creo un método que calcula la similitud coseno entre dos vectores y lo llamo tantas veces como documentos tenga. Me queda un array de valores, los scores, a mayor score, más parecido es ese documento a la query.

Ordeno con un Array de Tuplas los documentos por valor de score y luego creo un array de SearchItem del mismo tamaño que la cantidad de documentos. En cada posición guardo el título, texto sin normalizar y score del documento en orden de importancia descendente. Esto provoca que de un documento  $n$  en adelante los siguientes tengan score 0, y ocupen espacio por gusto, pero eso estoy intentando solucionarlo.

Al final SearchResult con mi array de SearchItem como argumento imprime en pantalla el Título y el Texto de los documentos ordenados de más a menos por valor de score.