# PROYECTO FINAL: SIMILITUD COSENO

## RECUPERACIÓN DE INFORMACIÓN
### BEATRIZ BELTRÁN MARTÍNEZ
### LICENCIATURA EN CIENCIAS DE LA COMPUTACIÓN
### PRIMAVERA 2022

GUSTAVO IVÁN MOLINA REBOLLEDO

000000000

Benemérita Universidad Autónoma de Puebla

ABSTRACT. In this paper, we present a Haskell program that calculates the similarity between every pair of given documents. The similarity is calculated using the cosine similarity. The documents are given as text files which are parsed and stored in a data structure. Then the program calculates the number of ocurrences of each word in the documents and the respective weight for each word. All of this is done in Haskell except the stemming process which is done through the use of a C library. We also discuss the performance of the program and the posible improvements that could be made.

3 de mayo de 2022

## I. Introducción

Este trabajo retoma las ideas que se han presentado en los trabajos anteriores. Sin embargo, se está tomando en cuenta una nueva organización del código, de forma que se espera una mejora con respecto a las entregas previas. Parte de estas mejoras es gracias a que se tiene una mejor visión de lo que se está realizando.

Al igual que en los documentos ya presentados, este trabajo se encuentra codificado en el lenguaje de programación Haskell. Se han hecho uso de varias librerías para la implementación de las funciones, mismas que se encuentran disponibles en Hackage. Esta elección tine un punto negativo que será explicado más adelante.

## II. Diseño del proyecto

El programa se encuentra dividido en varias secciones internas: *Parsing*, *Stage0*, *Stage1*, *Stage2*. Cada una de estás secciones se encargan de llevar a cabo trabajos específicos del procesado de los documentos. Además, se encuentra un modulo *main* que se encarga de la ejecución del programa.

### 2.1. Parsing

Para el *Parsing* asumimos que el documento se encuentra previamente almacenado en una estructura de datos de tipo *Text*. Ya que este proceso es de naturaleza impura, no se puede realizar en esta sección; *main* se encargará de ello.

Nuestro punto de entrada para esta sección es la función *parseDocument*. Esta función se encarga de extraer los datos del archivo de texto y de preprocesarlos para que puedan ser utilizados por el resto de las secciones.

El primer paso es hacer uso de la libraria *Text.XML.Light* para leer el archivo XML en una estructura manejable. El problema de hacer esto directamente es que la librería toma en cuenta los saltos de línea al momento de leer las etiquetas. Por lo tanto, se debe realizar un filtro para eliminar todos aquellos objetos que no sean elementos con información relevante.

Una vez que tenemos está divisíon de datos, se procede a dividir los datos en dos partes: la información del documento y el texto de la noticia. Esto se retorna en una lista de tuplas que puede

ser accedida por los siguientes pasos.

Lo siguiente a hacer es limpiar y preprocesar el texto de la noticia. Para ello eliminamos todos los números de la noticia, los signos de puntuación, las urls y se convierte el texto a minúsculas. Posteriormente, el texto es dividido en palabras y se filtra de forma que no se incluyan *stopwords*. Por último, se ensambla el texto en un solo *Text* y se da como entrada al algoritmo de *stemming*.

El algoritmo de *stemming* que estamos usando en este trabajo es un proyecto antiguo, además de que se trata de una biblioteca que hace *ffi* a una biblioteca en C. Una de las complicaciones al usar este algoritmo es que no se puede hacer uso de la librería de forma directa, sino que hay que modificar el código de la biblioteca para que funcione en un compilador de Haskell más moderno.

El algoritmo usado, a pesar de que es una biblioteca de C «antigua», compila sin problemas en Darwin moderno.

## 2.2. *Stage0 y representación interna*

Lo más interesante (e importante) de esta sección es que se ha diseñado un tipo producto para la representación de los datos que nos permite tener un mejor control y acceso en las etapas siguientes. Esta estructura está definida como se muestra a continuación:

```
data Documents = Documents {
  __docsN :: Int,
  __docs :: M.Map Int [(Int, Int)],
  __weights :: M.Map Int [(Int, Float)],
  __words :: [T.Text]
} deriving (Show, Generic)
```

disponibles (*__docsN*) y la lista de las palabras que se usan a lo largo de todos los documentos (*__words*). En cuanto a *__docs*, se define como una representación de indice invertido, pero manejando un *Map* para tener orden en el acceso. Está ultima guarda las apariciones de cada palabra en cada documento (si el documento la contiene).

Tener *__docs* y *__weights* puede parecer un poco ineficiente. Técnicamente lo es, pero en este

proyecto se requiere para poder mostrar los datos calculados por el programa. Además, en un principio requerimos de __*docs* para poder tener suficiente información para el calculo de __*weights*. Por lo que pensamos que su uso (en este caso específicamente) es justificado.

Además se definen una serie de *lenses* y morfismos para tener un acceso más fácil a los datos de la estructura. Esto es una mejora muy grande con respecto a las versiones anteriores que se han realizado. Simplifica grandemente la claridad del código y reduce la necesidad de declarar más funciones.

En cuanto a la generación, en esta etapa sólo se inicializa la estructura con el número de documentos y la lista de palabras en orden. Además se establecen los mapas con listas vacías.

### 2.3. Stage1

En esta sección no hay necesidad de una explicación muy extensa, ya que se trata de una etapa que se encarga de llenar la estructura con el número de apariciones de cada palabra en cada documento. Simplemente suministramos el texto de cada documento y la estructura de datos que se encuentra en la etapa anterior. El resultado de esta etapa es una estructura de datos con la información actualizada.

### 2.4. Stage2

Aquí nos ecargamos de calcular el peso de cada palabra en cada documento. En este caso basta con proporcionar sólo la estructura de datos de la etapa anterior, ya que contiene la información necesaria para calcular cada uno de los pesos.

#### 2.4..1.  TF

Definimos la función *tf* que calcula la frecuencia de una palabra en un documento de la siguiente forma:

```
idf :: Int -> Documents -> Float
idf word doc = let
  n = fromIntegral $ doc ^. _docsN
  d = fromIntegral $ length $ doc ^. _docs &
```

```
   M.lookup word & fromMaybe []
 in logBase 2 $ n / d
```

### 2.4..2. IDF

Definimos la función *idf* que calcula la frecuencia de una palabra en todos los documentos de la siguiente forma:

```
idf :: Int -> Documents -> Float
idf word doc = let
  n = fromIntegral $ doc ^. _docsN
  d = fromIntegral $ length $ doc ^. _docs &
    M.lookup word & fromMaybe []
  in logBase 2 $ n / d
```

### 2.4..3. TF-IDF

Finalmente, definimos la función *tfidf* que calcula el pesado del termino de la siguiente forma:

```
tfidf :: Int -> Int -> Documents -> Float
tfidf word doc docs =
  tf word doc docs * idf word docs
```

### 2.4..4. Peso

Simplemente basta con proveer de los indice necesarios (y de la estructura de datos de la etapa anterior) para calcular el peso de cada palabra. Esto se puede optener con facilidad usando los lenses. Además de que la estructura de la función *tfidf* es más similar a su definición teoríca.

Al acabar el procesor para todos los indices, tenemos una estructura actualizada con los valores para _ *weights* correspondientes a cada palabra.

## 2.5. Similitud Coseno

La similitud coseno es una de las más importantes funciones que se presentan en este proyecto, ya que es la que se encarga de calcular la similitud entre dos documentos. Pero en este caso, se calculá todos los pares de documentos. Esta función se define como sigue:

```
cosineSimilarity :: [Float] -> [Float] -> Float
cosineSimilarity v1 v2 = let
 t = sum $ zipWith (*) v1 v2
 a b = sqrt $ sum $ map (^2) b in
 t / (a v1) * (a v2)
```

Sólo hay que suministrarle los respectivos vectores de pesos de cada documento a comparar.

Sin embargo, la función no se ejecuta directamente, sino que es invocada por una función que se encarga de ejecutarla en cada par de documentos. Esto se realiza facilmente con el uso de una *comprehension list* que da los resultados de todos los cálculos.

## 2.6. Ensamblaje (main)

La longitud de los datos es demasiado grande para que se pueda calcular directamente a mano. Por lo que se necesita una función que se encarga de ejecutar todos los procesos necesarios para calcular la similitud entre todos los documentos.

En este apartado nos encontramos con funciones que calculan los datos de cada etapa y que se encargan de escribir los datos en un archivo compilable por LaTeX. Esto se realiza con el uso del paquete *HaTeX*, que nos permite escribir archivos de texto en LaTeX directamente desde el código en Haskell.

Por lo que la impresión de los datos en las siguientes secciones es hecha programáticamente.

## III. Tabla del conteo de documentos y palabras

## IV. Tabla del peso de documentos y palabras

This page contains a large numerical data matrix (cosine similarity values) that is too small and dense to transcribe reliably. The leftmost column contains word labels and the remaining columns contain numerical values (mostly 0.0).

## V. Tabla de la Similitud Coseno

## VI. Ejecución del programa

El programa se ejecuta desde la línea de comandos haciendo uso de *ghci*, o bien, compilando el programa con *ghc*. El resultado de la compilación nos da un programa que es capaz de generar cada uno de los pasos y de escribir los datos en un archivo LATEX compilable.

Para más información, consulte el archivo *README.md*

## VII. Conclusiones

Algo importante que se debe tener en cuenta es que el proceso de indexación es muy lento por la forma en la que se realiza el cálculo de los textos. Por lo que es recomendable que en una implementación más eficiente, se haga uso de estructuras que permitan el acceso de datos más rápido.

Haskell es un lenguaje de programación inmutable, por lo que no se puede modificar el contenido de una estructura de datos sin crear una nueva. Este es un problema ya que nos obliga a crear una nueva estructura para cada cambio que se realice.

Obviamente estos cambios afectan considerablemente al rendimiento de la aplicación. Sin embargo, simepre hay alternativas que podrían resolver este problema, como el uso de monadas, IO, estructuras ligadas a tipos lineales, etc.

Pero todas estas alternativas son muy complicadas y no son parte del alcance de este proyecto.

Otro punto a tomar en cuenta es que si bien el proceso de similitud debería ser muy rápido, esto no es lo que se ha observado. En realidad el cálculo de la similitud es bastante lento, siendo la parte más lenta de toda la aplicación.

Esto se debe probablemente a que hay que generar una gran cantidad de vectores. Un total de $101 * 101$ resultados son generados, lo que es bastante grande y puede ser complicado de calcular.

Estos tiempos son aún más lentos cuando se ejecuta desde el comando *ghci*. Por lo que es recomendable que se compile el programa con *ghc* para que se ejecute más rápido la generación de los documentos.