

## **Laboratorio 5: Representación tf y tf-idf**

Huerta Aguilar, Jesus., Huitzil Juárez, Guadalupe Quetzalli., Pérez Sánchez, Jasmine.,

Ruiz Ramírez, Gabino

Benemérita Universidad Autónoma de Puebla

Recuperación de la Información

Fecha: 27 Febrero 2024

**Resumen:** Con ayuda de los resultados obtenidos en los laboratorios 1,2,3 y 4 realizaremos una matriz mostrando la frecuencia de las palabras que hay en cada documento, para ello nos apoyaremos del vocabulario obtenido en el laboratorio 4, una vez obtenido el tf, (frecuencia de termino), calcularemos el idf-tf, para ello primo con ayuda de nuestra formula aprendida en clase calcularemos el idf y por último lo multiplicaremos con el tf.

**Palabras clave:** tokens, palabras vacías, re, nltk, Python, porter, lower, truncar, Split, stopwords, NPL, vocabulario, sorted, tf, idf-tf.

### **I. INTRODUCCIÓN**

En esta práctica, exploraremos en detalle dos aspectos fundamentales de la recuperación de información: la frecuencia de término (TF) y la frecuencia de término-inversa de documento (TF-IDF). Estos conceptos proporcionan una base sólida para comprender cómo se representa y se mide la importancia de los términos en un conjunto de documentos. Recordemos que el TF mide la frecuencia con la que un término específico aparece en un documento, mientras que el TF-IDF, combina el TF con la Inversa de Frecuencia de Documento (IDF), que mide la rareza de un término en el corpus de documentos. La multiplicación del TF y el IDF permite dar más peso a los términos que son frecuentes en un documento específico pero raros en el corpus en general, lo que ayuda a identificar términos más relevantes.

### **II. OBJETIVO Y PLANTEAMIENTO DEL PROBLEMA**

El objetivo principal de esta práctica es familiarizarse con el modelo vectorial en la recuperación de información y comprender cómo se representan documentos y consultas utilizando este enfoque. Además, se busca aplicar los conceptos de Frecuencia de Término (TF) e Inversa de Frecuencia de Término en Documento (IDF-TF) para mejorar la representación y la relevancia de los términos en el proceso de recuperación de información.

### III. DESARROLLO EXPERIMENTAL

Para realizar esta práctica necesitaremos:

- Software Python
- Librería NLTK
- Colección NPL de la página: [https://ir.dcs.gla.ac.uk/resources/test\\_collections/npl/](https://ir.dcs.gla.ac.uk/resources/test_collections/npl/)
- Laboratorios 1, 2, 3 y 4.

Para poder realizar nuestro código debemos tomar en cuenta dos cosas:

- La dimensión de los vectores será igual al tamaño del vocabulario reducido que obtuvimos en la práctica 4.
- Nos apoyaremos del documento en que se encuentra el vocabulario reducido.

Contemplando estos dos puntos procederemos a desarrollar nuestro código.

1. Primero representaremos cada documento empleando el modelo vectorial y el esquema de pesado tf y guardamos la matriz generada.

Creamos nuestra función que permite leer el contenido del documento, en seguida creamos una función que permita generar la matriz en un documento de acuerdo con el tamaño del vocabulario, tomando en cuenta la línea del documento.

En seguida recuperamos el vocabulario guardado en el archivo de vocabulario, este lo guardaremos en una lista para que sea más fácil al momento de emplearlo.

```
15 # leer un archivo y devolver su contenido como una lista de líneas
16 def leer_archivo(nombre_archivo):
17     with open(nombre_archivo, 'r', encoding='utf-8') as archivo:
18         contenido = archivo.readlines()
19     return contenido
20
21 # escribir una matriz en un archivo de texto
22 def escribir_matriz(nombre_archivo, matriz):
23     with open(nombre_archivo, 'w', encoding='utf-8') as archivo:
24         contador = 0
25         for fila in matriz:
26             contador += 1
27             progress_bar(contador, 11419) # LLAMADA A FUNCIÓN 'progress_bar()'
28             archivo.write('\t'.join(map(str, fila)) + '\n')
29
30 # generar el vocabulario a partir de un archivo de texto
31 def generar_vocabulario(archivo_vocabulario):
32     vocabulario = {}
33     lineas = leer_archivo(archivo_vocabulario)
34     contador = 0
35     for linea in lineas:
36         contador += 1
37         palabras = linea.strip().split()
38         for palabra in palabras:
39             vocabulario[palabra] = len(vocabulario)
40     return vocabulario
```

En seguida vectorizaremos el documento es decir asignaremos a cada documento un vector numérico que capture alguna característica relevante del texto. Para ello creamos una función que permita realizar este proceso, donde crearemos una matriz para almacenar las representaciones vectoriales, creamos una lista de frecuencias de cada palabra del vocabulario que se encuentre en el documento con ayuda de un contador, una vez obtenida la frecuencia de la palabra en el documento actual este valor se guardara en la matriz.

```

46 def vectorizar_documentos(archivo_documentos, vocabulario):
47     lineas = leer_archivo(archivo_documentos)
48     matriz = []
49     print("\n\nVECTORIZANDO DOCUMENTOS...")
50     for linea in lineas:
51         frecuencias = [0] * len(vocabulario)
52         palabras = linea.strip().split()
53         for palabra in palabras:
54             if palabra in vocabulario:
55                 indice = vocabulario[palabra]
56                 frecuencias[indice] += 1
57         matriz.append(frecuencias)
58     return matriz

```

Una vez creada nuestras funciones a ocupar procederemos a guardar nuestros documentos a ocupar en variables para hacer más fácil el manejo de estos, por ejemplo, el archivo del vocabulario, el archivo donde se encuentran los documentos, y el archivo donde se guardará la matriz dada. En seguida llamaremos a la función generar vocabulario creada anteriormente, vectorizaremos los documentos y por último, guardaremos la matriz en el archivo de salida llamado matriz-docs.txt. Para corroborar que la matriz cuente con el numero correcto de filas y columnas imprimiremos ese valor para compararlo después con los documentos.

```

64 # Archivos de entrada
65 archivo_vocabulario = r'C:\Users\mini_OneDrive\Documents\Code Test\TEST 1\lab5\vocabularioReducidoT.txt'
66 archivo_documentos = r'C:\Users\mini_OneDrive\Documents\Code Test\TEST 1\lab5\documento-TRUNCADO.txt'
67 # Archivo de salida
68 archivo_salida = r'C:\Users\mini_OneDrive\Documents\Code Test\TEST 1\lab5\matriz-docs.txt'
69
70 # Generar vocabulario
71 vocabulario = generar_vocabulario(archivo_vocabulario)
72
73 # Vectorizar documentos
74 matriz = vectorizar_documentos(archivo_documentos, vocabulario)
75
76 # Escribir matriz en archivo de salida
77 print("\n\nCREANDO MATRIZ TF...")
78 escribir_matriz(archivo_salida, matriz)
79
80 print("\n\nMATRIZ TF GENERADA Y GUARDADA EN:", archivo_salida)
81 # Imprimir la cantidad de filas y columnas de la matriz
82 print("\n\nFILAS DE LA MATRIZ TF:", len(matriz))
83 print("\n\nCOLUMNAS DE LA MATRIZ TF:", len(matriz[0]))
84

```

2. De la misma manera que el paso anterior, representaremos ahora el archivo de queries empleando el modelo vectorial y esquema de pesado tf.

Cambiando únicamente el valor del total de los documentos(filas) pasando de 11419 a 93

```

21 # escribir una matriz en un archivo de texto
22 def escribir_matriz(nombre_archivo, matriz):
23     with open(nombre_archivo, 'w', encoding='utf-8') as archivo:
24         contador = 0
25         for fila in matriz:
26             contador += 1
27             progress_bar(contador, 93) # LLAMADA A FUNCIÓN 'progress_bar()'
28             archivo.write('\t'.join(map(str, fila)) + '\n')

```

También cambiando el valor donde guardaremos el resultado en un nuevo archivo llamado matriz-query.txt que en variable también es archivo de salida.

```

64 # Archivos de entrada
65 archivo_vocabulario = r'C:\Users\mini_OneDrive\Documentos\Code Test\TEST 1\lab5\vocabularioReducidoT.txt'
66 archivo_documentos = r'C:\Users\mini_OneDrive\Documentos\Code Test\TEST 1\lab5\cueri-TRUNCADO.txt'
67 # Archivo de salida
68 archivo_salida = r'C:\Users\mini_OneDrive\Documentos\Code Test\TEST 1\lab5\matriz-query.txt'
69

```

3. Repita los pasos 1 y 2 con esquema de pesado tf-idf para documento y para queries. Para este caso crearemos una función que permita calcular el IDF, primero calculando el número de documentos (N) que hay en la colección, el número de palabras y el número de documentos con esa palabra ( $N_t$ ), aplicamos la formula aprendida en clase para calcular el idf:

$$idf = \log \frac{N}{N_t} + 1$$

Después crearemos una función que permita escribir el vocabulario junto con su idf y esto lo guardaremos en un archivo llamado IDF-docs, en variable llamado `archivo_idf`

```

89 # Función para calcular el IDF
90 def calcular_idf(matriz):
91     num_documentos = len(matriz)
92     num_palabras = len(matriz[0])
93     idf = [0] * num_palabras
94     print("\nCALCULANDO IDF...")
95     for i in range(num_palabras):
96         progress_bar(i + 1, num_palabras) # LLAMADA A FUNCIÓN 'progress_bar()'
97         # Contar en cuántos documentos aparece la palabra
98         num_documentos_con_palabra = sum([1 for documento in matriz if documento[i] > 0])
99         # Calcular el IDF
100         idf[i] = math.log(num_documentos / (num_documentos_con_palabra)) + 1
101     return idf
102
103 # Calcular IDF
104 idf = calcular_idf(matriz)
105
106 # Función para escribir el vocabulario junto con sus valores IDF en un archivo de texto
107 def escribir_vocabulario_idf(nombre_archivo, vocabulario, idf):
108     with open(nombre_archivo, 'w', encoding='utf-8') as archivo:
109         for palabra, valor_idf in zip(vocabulario, idf):
110             archivo.write(palabra + '\t' + str(valor_idf) + '\n')
111
112 # Generar archivo con vocabulario e IDF
113 archivo_idf = r'C:\Users\mini_OneDrive\Documentos\Code Test\TEST 1\lab5\IDF-docs.txt'
114 escribir_vocabulario_idf(archivo_idf, vocabulario, idf)
115
116 print("\n\nVALORES IDF GUARDADOS EN:", archivo_idf)

```

Una vez obtenido el idf procederemos a multiplicarlo por el tf para así obtener el IDF-TF guardando la matriz obtenida en el archivo de salida llamado `matriz-TF-IDF-docs.txt`.

```

122 # Función para calcular el TF-IDF
123 def calcular_tfidf(matriz_tf, idf):
124     print("\nCALCULANDO TF - IDF...")
125     matriz_tfidf = []
126     contador = 0
127     for tf_vector in matriz_tf:
128         contador += 1
129         tfidf_vector = [tf * idf[i] for i, tf in enumerate(tf_vector)]
130         matriz_tfidf.append(tfidf_vector)
131         progress_bar(contador, 11419) # LLAMADA A FUNCIÓN 'progress_bar()'
132     return matriz_tfidf
133
134 # Calcular TF-IDF
135 matriz_tfidf = calcular_tfidf(matriz, idf)
136
137 # Archivo de salida para la matriz TF-IDF
138 archivo_salida_tfidf = r'C:\Users\mini_OneDrive\Documentos\Code Test\TEST 1\lab5\matriz-TF-IDF-docs.txt'
139
140 # Escribir matriz TF-IDF en archivo de salida
141 print("\n\nCREANDO MATRIZ TF-IDF...")
142 escribir_matriz(archivo_salida_tfidf, matriz_tfidf)
143
144 print("\n\nMATRIZ TF-IDF GENERADA Y GUARDADA EN:", archivo_salida_tfidf)

```

De la misma manera calcularemos el TF-IDF para los queries: únicamente cambiando los archivos en los que se guardará:

```

115 # Generar archivo con vocabularioclids IDF
116 archivo_idf = r'C:\Users\mini_OneDrive\Documentos\Code Test\TEST 1\lab5\IDF-query.txt'
117 escribir_vocabulario_idf(archivo_idf, vocabulario, idf)
118
119 print("\n\nVALORES IDF GUARDADOS EN:", archivo_idf)
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140 # Archivo de salida para la matriz TF-IDF
141 archivo_salida_tfidf = r'C:\Users\mini_OneDrive\Documentos\Code Test\TEST 1\lab5\matriz-TFIDF-query.txt'
142
143 # Escribir matriz TF-IDF en archivo de salida
144 print("\n\nCREANDO MATRIZ TF-IDF...")
145 escribir_matriz(archivo_salida_tfidf, matriz_tfidf)
146
147 print("\n\nMATRIZ TF-IDF GENERADA Y GUARDADA EN:", archivo_salida_tfidf)

```

## IV. DISCUSIÓN Y RESULTADOS

Como resultados de nuestra primera matriz de TF obtuvimos lo siguiente:

Para documentos al momento de ejecutar obtenemos que el número de filas y columnas es correcto, ya que consta de 11429 documentos y 7616 términos:

```
C:\Users\mini_>C:\Users\mini\AppData\Local\Programs\Python\Python312\python.exe "C:\Users\mini_OneDrive\Documents\Code Test\TEST 1\lab5\lab5-docs.py"

VECTORIZANDO DOCUMENTOS...

CREANDO MATRIZ TF...
| 100.00%

MATRIZ TF GENERADA Y GUARDADA EN: C:\Users\mini_OneDrive\Documents\Code Test\TEST 1\lab5\matriz-docs.txt

FILAS DE LA MATRIZ TF: 11429
COLUMNAS DE LA MATRIZ TF: 7616

CALCULANDO IDF...
| 100.00%

VALORES IDF GUARDADOS EN: C:\Users\mini_OneDrive\Documents\Code Test\TEST 1\lab5\IDF-docs.txt

CALCULANDO TF - IDF...
| 100.00%

CREANDO MATRIZ TF-IDF...
| 100.00%

MATRIZ TF-IDF GENERADA Y GUARDADA EN: C:\Users\mini_OneDrive\Documents\Code Test\TEST 1\lab5\matriz-TFIDF-docs.txt

C:\Users\mini_>
```

Nuestra matriz es la siguiente observamos que la mayor parte está conformada por 0 ya que la mayor parte del vocabulario no se encuentra en cada documento de esta forma observamos que es más difícil encontrar un 1 u otro valor:

Para queries al momento de ejecutar podemos observar de igual manera que el valor es correcto donde las filas es el número de consulta mientras que las columnas corresponden al número de termino:

```
Command Prompt
C:\Users\mini_>C:\Users\mini\AppData\Local\Programs\Python\Python312\python.exe "C:\Users\mini_OneDrive\Documentos\Code Test\TEST 1\lab5\lab5-query.py"
VECTORIZANDO DOCUMENTOS...

CREANDO MATRIZ TF...
| 100.00%
MATRIZ TF GENERADA Y GUARDADA EN: C:\Users\mini_OneDrive\Documentos\Code Test\TEST 1\lab5\matriz-query.txt
FILAS DE LA MATRIZ TF: 93
COLUMNAS DE LA MATRIZ TF: 7616

CALCULANDO IDF...
| 100.00%
VALORES IDF GUARDADOS EN: C:\Users\mini_OneDrive\Documentos\Code Test\TEST 1\lab5\idf-query.txt
CALCULANDO TF - IDF...
| 100.00%
CREANDO MATRIZ TF-IDF...
| 100.00%
MATRIZ TF-IDF GENERADA Y GUARDADA EN: C:\Users\mini_OneDrive\Documentos\Code Test\TEST 1\lab5\matriz-TFIDF-query.txt
C:\Users\mini_>
```

Así mismo observamos de igual manera que la matriz esta conformada por mas 0 que 1 u otro valor:



TEST 1	lab5	matriz-query.txt
1	0	0
2	0	0
3	0	0
4	0	0
5	0	0
6	0	0
7	0	0
8	0	0
9	0	0
10	0	0
11	0	0
12	0	0
13	0	0
14	0	0
15	0	0
16	0	0
17	0	0
18	0	0
19	0	0
20	0	0
21	0	0
22	0	0
23	0	0
24	0	0
25	0	0
26	0	0
27	0	0
28	0	0
29	0	0
30	0	0
31	0	0
32	0	0

Como IDF obtuvimos el siguiente:

Documentos aquí podemos observar el término junto con su IDF, para corroborar el resultado tomaremos en cuenta el termino 23 con su IDF:

The screenshot shows a code editor with a file explorer on the left. The file explorer shows a project structure with folders like 'TEST 1' and files like 'lab5-docs.py', 'lab5-query.py', and 'lab5-queryado.py'. The main editor displays a list of terms and their IDF values. The term 'absent' is highlighted in blue.

Term	IDF
aachen	10.343909263897597
aad	9.650762083337654
aan	9.650762083337654
abac	9.245296975229488
abandon	10.343909263897597
abbrevi	9.245296975229488
abd	10.343909263897597
abel	9.650762083337654
aber	10.343909263897597
aberr	8.397999114842285
abil	8.397999114842285
abing	9.245296975229488
abrid	10.343909263897597
abridg	10.343909263897597
abrupt	8.734471351463498
abruptii	8.957614982777786
abscessa	9.650762083337654
absenc	6.705323184171212
absent	8.552149794669543
absolut	7.011704753722394
absorb	6.216774878852506
absorpt	4.588167058310680
abstract	6.472708252989787
abu	9.650762083337654
abund	8.552149794669543
academi	8.734471351463498

Queries: de igual forma que en los documentos podemos observar el término con su IDF correspondiente y para corroborar tomaremos en cuenta el IDF del término resaltado en azul y lo buscaremos en el siguiente paso

The screenshot shows a code editor with a file explorer on the left. The file explorer shows a project structure with folders like 'TEST 1' and files like 'lab5-docs.py', 'lab5-query.py', and 'lab5-queryado.py'. The main editor displays a list of terms and their IDF values. The term 'absent' is highlighted in blue.

Term	IDF
aachen	1
aad	1
aan	1
abac	1
abandon	1
abbrevi	1
abd	1
abel	1
aber	1
aberr	1
abil	1
abing	1
abrid	1
abridg	1
abrupt	1
abruptii	1
abscessa	1
absenc	1
absent	1
absolut	1
absorb	1
absorpt	5.532598493150156
abstract	4.839452312593311
abu	1
abund	1
academi	1

Por último, nuestro resultado de matriz TF-IDF fue el siguiente:

Para documentos recordando el valor anterior y ajustando la matriz podemos observar que en el documento 2732 (fila) se encuentra el resultado TF-IDF correspondiente al termino(columna):



Para queries: recordando el término con su IDF seleccionado anteriormente podemos observar que en la consulta 21 se encuentra el resultado correcto y correspondiente al término; en este caso era el termino 28 y ajustando la matriz podemos observar que se encuentra correctamente.

## V. CONCLUSIONES

En resumen, en este laboratorio hemos explorado la representación de documentos y consultas en el modelo vectorial, utilizando herramientas como Python y NLTK. A partir de los archivos pre-procesados de la colección NPL y los queries proporcionados, hemos aplicado el esquema de pesado tf y tf-idf para crear matrices que representan tanto nuestros documentos como nuestras consultas.



Al emplear el esquema de pesado tf, pudimos asignar pesos a cada término en función de su frecuencia en los documentos, lo que nos brinda una manera efectiva de representar la importancia relativa de cada palabra dentro de un texto. Por otro lado, al aplicar el esquema tf-idf, consideramos tanto la frecuencia del término en el documento como su importancia en toda la colección, lo que nos permite capturar mejor la relevancia de un término en un contexto más amplio.

Estas representaciones vectoriales nos brindan una manera eficaz de comparar documentos y consultas en un espacio matemático, facilitando la implementación de algoritmos de recuperación de información. Al guardar las matrices generadas, hemos creado una base de datos eficiente que podemos utilizar en futuras prácticas y análisis.

Este laboratorio ha sido fundamental para comprender cómo las técnicas de modelado de texto nos permiten estructurar y analizar grandes colecciones de documentos.

## VI. BIBLIOGRAFÍA

1. NLTK. (2022). SnowballStemmer. Recuperado de <https://www.nltk.org/api/nltk.stem.SnowballStemmer.html?highlight=stopwords>
2. NLTK. (2022). WordNetLemmatizer. Recuperado de <https://www.nltk.org/api/nltk.stem.WordNetLemmatizer.html?highlight=wordnet>