

Certamen III: Lenguajes de Programación

Implementación de un Motor de Búsqueda mediante Índice Invertido y Procesamiento Recursivo

Camilo Chavol – Graciela Suárez - Romily Barría

Prof. Guía: Alonso Inostrosa Psijas

Escuela de Ingeniería Informática

Facultad de Ingeniería

Universidad de Valparaíso

Resumen

Este proyecto consiste en diseñar y construir un Motor de Búsqueda de texto completo funcional, partiendo desde los datos crudos (dataset.csv) hasta la consulta interactiva. Para ello, implementamos la estructura de datos más eficiente para la recuperación de información: el Índice Invertido.

- El proyecto exigió la integración exitosa de múltiples paradigmas, destacando la programación declarativa.
- Modularidad: El flujo de trabajo se dividió en tres módulos (Indexación, Limpieza, Búsqueda) para garantizar la escalabilidad y el cumplimiento de requisitos de lenguaje específicos.
 - Recursividad: Se implementó obligatoriamente el filtro de stopwords mediante un algoritmo recursivo en Python, demostrando el manejo del flujo de control funcional.
 - Precisión: La búsqueda se basa en la Lógica AND, utilizando operaciones de conjuntos (Intersección) para asegurar que el resultado coincida con todos los términos ingresados por el usuario.

Índice

- Introducción
- Bases y Supuestos
- Arquitectura del Sistema
- Implementación Recursiva
- Lógica de Búsqueda
- Resultados
- Conclusión



Introducción

Introducción

El objetivo principal es resolver el problema de la recuperación de información eficiente. En lugar de recorrer todos los textos secuencialmente para cada búsqueda (lo cual sería ineficiente), se pre-procesa la información en un Índice Invertido.

Conceptos Clave:

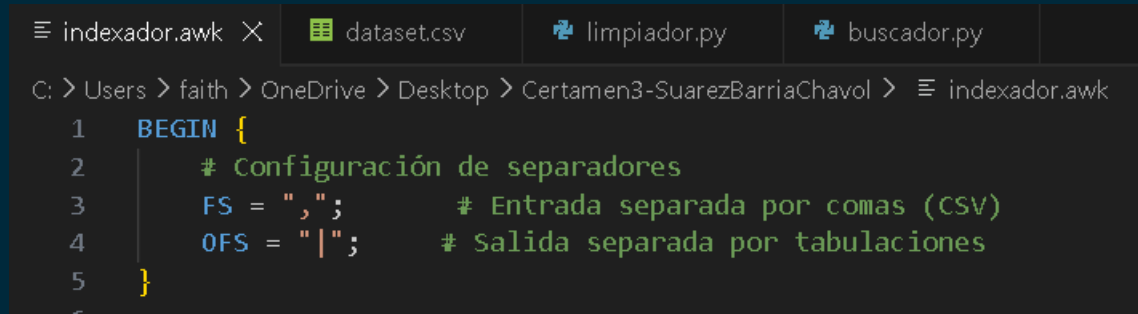
- Vocabulario (V): Conjunto de palabras únicas presentes en el corpus.
- Posting List: Lista de identificadores de documentos asociados a una palabra.
- Stopwords: Palabras como "el", "la", "y", que deben ser eliminadas para optimizar el almacenamiento y la relevancia.

La solución propuesta automatiza el flujo de trabajo (Pipeline) desde la lectura de datos crudos hasta la consulta del usuario final.

Bases y Supuestos

Bases y Supuestos

- Inicialización (Patrón BEGIN): Se utiliza el bloque BEGIN para definir variables antes de procesar la primera línea del archivo, garantizando la configuración global del entorno.
- Separador de Entrada (FS): Se define la variable incorporada FS (Field Separator) como coma (,), ya que el dataset de entrada es un archivo CSV estándar.
- Separador de Salida (OFS): Se establece la variable OFS (Output Field Separator) como una barra vertical (|). Esto facilita el procesamiento posterior en Python, evitando conflictos con comas dentro del texto.

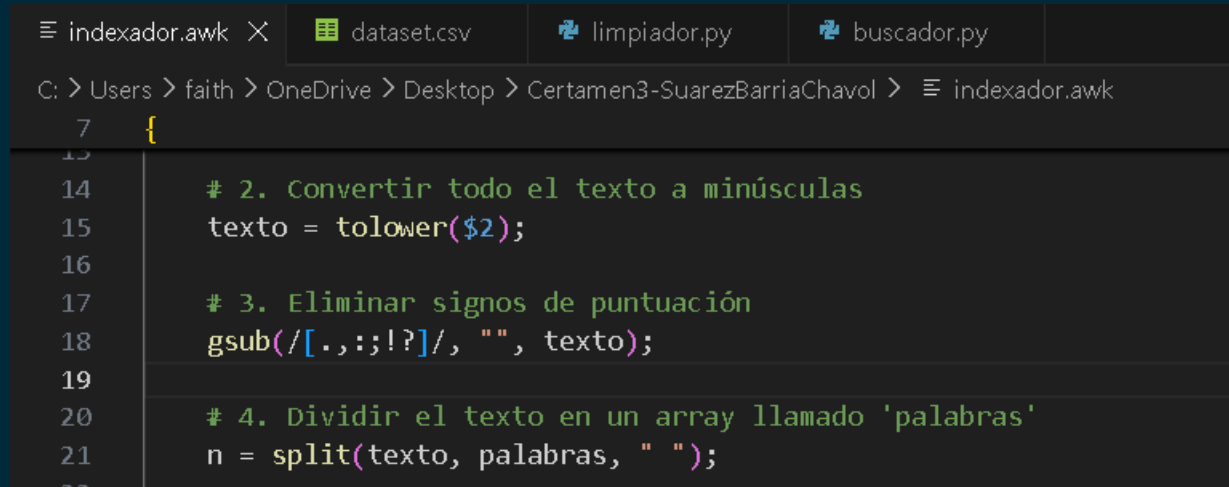


The screenshot shows a code editor with four tabs: 'indexador.awk', 'dataset.csv', 'limpiador.py', and 'buscador.py'. The active tab is 'indexador.awk', which contains the following code:

```
C: > Users > faith > OneDrive > Desktop > Certamen3-SuarezBarriaChavol > indexador.awk
1 BEGIN {
2     # Configuración de separadores
3     FS = ",";      # Entrada separada por comas (CSV)
4     OFS = "|";     # Salida separada por tabulaciones
5 }
```

Bases y Supuestos

- Estandarización (Case-Folding): Se utiliza la función `tolower` para convertir todo el texto a minúsculas. Esto elimina la sensibilidad a mayúsculas (case-sensitivity), asegurando que "Perro" y "perro" se traten como el mismo término.
- Limpieza de Ruido: Mediante la función `gsub` y expresiones regulares, se eliminan los signos de puntuación (.,:;!?), dejando únicamente los caracteres alfanuméricos relevantes para la indexación.
- Tokenización (Split): La función `split` divide el texto limpio en un arreglo de palabras (palabras), permitiendo procesar cada término individualmente dentro del bucle.



```
indexador.awk X dataset.csv limpiador.py buscador.py
C: > Users > faith > OneDrive > Desktop > Certamen3-SuarezBarriaChavol > indexador.awk
7 {
14 # 2. Convertir todo el texto a minúsculas
15 texto = tolower($2);
16
17 # 3. Eliminar signos de puntuación
18 gsub(/[.,:;!?]/, "", texto);
19
20 # 4. Dividir el texto en un array llamado 'palabras'
21 n = split(texto, palabras, " ");
22 }
```




Arquitectura del Sistema

Arquitectura del Sistema

-Flujo de Trabajo Modular: El sistema se diseña como una tubería (pipeline) de tres etapas secuenciales, desacoplando la extracción de datos de la lógica de búsqueda.

-Orquestación (Makefile): Se utiliza un archivo Makefile para automatizar la ejecución. Esto garantiza que la indexación (AWK) ocurra estrictamente antes de la limpieza y la búsqueda (Python).

-Interoperabilidad: La comunicación entre los módulos (AWK y Python) se realiza mediante archivos de texto plano (index_raw.txt), lo que permite integrar distintos paradigmas de programación de manera transparente.

```
makefile.txt X
C: > Users > faith > OneDrive > Desktop > Certamen3-SuarezBarriaChavol > makefile.txt
1  # Makefile para Certamen 3 - Lenguajes de Programación
2
3  run:
4      @echo "[PASO 1] Indexacion"
5      awk -f indexador.awk dataset.csv > index_raw.txt
6
7      @echo "\n[PASO 2] Eliminando stopwords (Python Recursivo)..."
8      python limpiador.py
9
10     @echo "\n[PASO 3] Iniciando Buscador..."
11     python buscador.py
12
13  clean:
14      # Limpia los archivos generados
15      del index_raw.txt index_master.txt
```



Implementación Recursiva

Implementación Recursiva

-Sustitución de Iteraciones: Se reemplazan los bucles tradicionales (for/while) por recursividad explícita, cumpliendo con el requisito de utilizar el paradigma funcional para controlar el flujo de ejecución.

-Caso Base: La función verifica si el índice actual ha alcanzado el total de líneas del archivo (index == len). Si se cumple, detiene la recursión y retorna el resultado acumulado.

-Paso Recursivo: Si no es el final, la función procesa la línea actual y se invoca a sí misma (return procesar...), incrementando el índice (index + 1) y pasando el estado actualizado como argumento, en lugar de modificar variables globales.

```
limpiador.py X buscador.py indexador.awk
C: > Users > faith > OneDrive > Desktop > Certamen3-SuarezBarriaChavol > limpiador.py > ...

25 def procesar_y_filtrar_recursivo(lineas_csv, stopwords, index=0, indice_filtrado=[]):
26     """
27     FUNCIÓN RECURSIVA (Requerimiento 4)
28     Procesa las líneas del CSV, tokeniza, y filtra stopwords.
29     """
30     # Caso Base
31     if index == len(lineas_csv):
32         return indice_filtrado
33
34     try:
35         doc_id = lineas_csv[index][0].strip()
36         texto = lineas_csv[index][1]
37
38         palabras = tokenizar_y_limpiar(texto)
39
40         for palabra in palabras:
41             if not palabra: continue
42
43             if palabra in stopwords:
44                 print(f"[Stopword Eliminada] '{palabra}' del documento {doc_id}")
45             else:
46                 # Agregamos la palabra limpia y el ID usando el separador '|'
47                 indice_filtrado.append(f"{palabra}|{doc_id}")
48
49     except IndexError:
50         pass # Ignoramos líneas mal formadas
51
52     # Llamada recursiva al siguiente índice
53     return procesar_y_filtrar_recursivo(lineas_csv, stopwords, index + 1, indice_filtrado)
54
```

Lógica de Búsqueda

Lógica de Búsqueda

-Persistencia en Memoria: Siguiendo los requerimientos, el índice se carga completamente en RAM al iniciar el programa, evitando lecturas a disco durante las consultas para maximizar la velocidad.

-Implementación del Índice Invertido: Se utiliza un Diccionario (Hash Map) como estructura principal.

- Clave (Key): Corresponde a cada palabra única del vocabulario (\$V\$).
- Valor (Value): Corresponde a la Lista de Posteo, implementada aquí como un Conjunto (set).

-Uso de Conjuntos (set): Se eligió esta estructura en lugar de listas simples para garantizar que no existan IDs de documentos duplicados y para optimizar matemáticamente la operación de intersección posterior.

```
limpiador.py  buscador.py  indexador.awk
C:\Users\faith> OneDrive\Desktopt> Certamen3-SuarezBarriaChavol > buscador.py > cargar_indice_invertido
3 def cargar_indice_invertido(ruta):
4     """
5     Carga el índice invertido en un Diccionario de Python.
6     Aplica limpieza defensiva al cargar las claves.
7     """
8     indice = {}
9     try:
10         # LECTURA: Usamos 'utf-8-sig' y el separador Pipe '|'
11         with open(ruta, 'r', encoding='utf-8-sig') as f:
12             for linea in f:
13                 partes = linea.strip().split('|')
14                 if len(partes) < 2: continue
15
16                 palabra = partes[0].strip()
17                 doc_id = partes[1].strip()
18
19                 if palabra not in indice:
20                     indice[palabra] = set()
21                 indice[palabra].add(doc_id)
22
23             return indice
24     except FileNotFoundError:
25         print("Error crítico: No se encontró el índice maestro. Asegúrese de ejecutar los pasos 1 y 2.")
26         sys.exit(1)
```

Lógica de Búsqueda

-Lógica Booleana (AND): El algoritmo implementa una búsqueda conjuntiva estricta. Un documento solo es considerado un resultado válido si contiene todos los términos ingresados en la consulta.

-Operación de Conjuntos:

- Se recupera el conjunto de documentos del primer término (indice.get).
- Se itera sobre los términos restantes, aplicando la operación matemática de Intersección (intersection) en cada paso.

-Reducción del Espacio de Búsqueda: Esta estrategia filtra progresivamente los candidatos, asegurando que el resultado final cumpla con la restricción de coincidencia total descrita en el problema.

```
limpiador.py  buscador.py  indexador.awk
C:\Users\faith> OneDrive\ Desktop\ Certamen3-SuarezBarriaChavol\ buscador.py > cargar_indice_invertido
27 def buscar(consulta, indice):
28     """
29     Realiza la búsqueda intersectando listas de posteo (Lógica AND).
30     """
31     palabras_consulta = consulta.lower().strip().split()
32
33     if not palabras_consulta:
34         return set()
35
36     primera_palabra = palabras_consulta[0]
37     documentos_resultado = indice.get(primera_palabra, set())
38
39     for palabra in palabras_consulta[1:]:
40         key_limpia = palabra.strip()
41         docs_siguiete = indice.get(key_limpia, set())
42         documentos_resultado = documentos_resultado.intersection(docs_siguiete)
43
44     return documentos_resultado
45
```

Lógica de Búsqueda

-Lógica Conjuntiva (AND): El motor implementa una búsqueda estricta. Un documento solo es relevante si contiene todos los términos de la consulta, descartando coincidencias parciales.

-Algoritmo de Intersección:

- Inicialización: Se recupera el conjunto de documentos asociado al primer término de la consulta (`indice.get`).
- Filtrado Progresivo: Se itera sobre los términos restantes, actualizando el conjunto de resultados mediante la operación matemática de Intersección (`intersection`).

-Eficiencia: Al reducir el espacio de búsqueda en cada paso, se garantiza que el resultado final cumpla con la restricción de coincidencia total sin necesidad de verificar documento por documento manualmente.

```
limpiador.py  buscador.py X  indexador.awk
C:\Users\faith> OneDrive\Desktop> Certamen3-SuarezBarriaChavol > buscador.py > cargar_indice_invertido
27 def buscar(consulta, indice):
28     """
29     Realiza la búsqueda intersectando listas de posteo (Lógica AND).
30     """
31     palabras_consulta = consulta.lower().strip().split()
32
33     if not palabras_consulta:
34         return set()
35
36     primera_palabra = palabras_consulta[0]
37     documentos_resultado = indice.get(primera_palabra, set())
38
39     for palabra in palabras_consulta[1:]:
40         key_limpia = palabra.strip()
41         docs_siguiente = indice.get(key_limpia, set())
42         documentos_resultado = documentos_resultado.intersection(docs_siguiente)
43
44     return documentos_resultado
45
```


Resultados

Resultados: Video de Funcionamiento

Resultado: Resultado Generado

-Interfaz de Usuario: El sistema permite el ingreso interactivo de consultas, procesando la entrada en tiempo real.

-Prueba de Un Término: La consulta "perro" recupera correctamente todos los documentos asociados, validando la integridad del índice invertido.

-Prueba de Intersección (AND):

- Consulta: "perro casa".
- Resultado: El sistema filtra correctamente y devuelve solo los documentos (ej. 1 y 5) que contienen ambos términos, demostrando la lógica de conjuntos estricta.

-Prueba de Múltiples Términos: La consulta "gato casa arbol" confirma la capacidad del algoritmo para reducir progresivamente el espacio de búsqueda hasta encontrar la coincidencia exacta.

```
PS C:\Users\faith\OneDrive\Desktop\Certamen3-SuarezBarriaChavo1> python buscador.py
>>> MOTOR DE BÚSQUEDA (ÍNDICE INVERTIDO) <<<
Índice cargado en memoria. 13 términos únicos.

Ingrese su búsqueda (o 'salir'): perro
RESULTADO: Documentos encontrados -> 1, 3, 5
-----
Ingrese su búsqueda (o 'salir'): perro casa
RESULTADO: Documentos encontrados -> 1, 5
-----
Ingrese su búsqueda (o 'salir'): gato casa arbol
RESULTADO: Documentos encontrados -> 2
-----
Ingrese su búsqueda (o 'salir'): salir

Fin del programa.
PS C:\Users\faith\OneDrive\Desktop\Certamen3-SuarezBarriaChavo1> █
```

Conclusión

Conclusión

-Integración de Paradigmas: El proyecto demuestra el trabajo en conjunto de distintos paradigmas: la eficiencia de Scripting (AWK) para el procesamiento de texto crudo (raw text, que son los datos de caracteres tal como se reciben de la fuente original) y la flexibilidad de la Programación Funcional (Python) para la lógica de negocio.

-Eficiencia del Índice Invertido: Se validó que esta estructura de datos es superior a la búsqueda secuencial, permitiendo tiempos de respuesta inmediatos gracias a la indexación previa y la carga en memoria.

-Cumplimiento de Requisitos:

- Modularidad: El diseño desacoplado (Pipeline) facilita el mantenimiento.
- Recursividad: Se implementó exitosamente el filtrado de stopwords mediante llamadas recursivas, cumpliendo con las exigencias del paradigma funcional.
- Precisión: La lógica de conjuntos garantizó resultados exactos bajo el modelo booleano AND.



**Gracias por su
Atención**