

UNIVERSIDAD CENTRAL DE VENEZUELA
FACULTAD DE CIENCIAS
POSTGRADO EN CIENCIAS DE LA COMPUTACIÓN



ANTEPROYECTO:

Solución Sistema Complementario Saber UCV (SSCSU): extracción de datos desde PDF, clasificación y construcción de un sistema de recuperación de información en línea

Anteproyecto presentado
por el Econ. José Miguel Avendaño Infante
Tutor: Dr. Andrés Sanoja

Caracas, Octubre de 2021

Resumen:

Se presenta la propuesta de un sistema denominado ****Solución Sistema Complementario Saber UCV (SSCSU)**** para hacer procesos de Recuperación de Información sobre los Resúmenes, tanto de las Tesis como de los Trabajos Especiales de Grado (TEG), que se encuentran alojados en el repositorio institucional Saber UCV (www.saber.ucv.ve).

Se aplican técnicas de Procesamiento de Lenguaje Natural (NLP), de Minería de Texto y de indexación en base de datos sobre los textos.

Esta propuesta se basa en que mediante la búsqueda de palabras o frases, aplicando filtros y determinando la granularidad, se genere un (*query*) con el que se puedan recuperar los trabajos en que se encuentran contenidas tales palabras y a partir de ahí enriquecer la experiencia del usuario con la presentación de la información recuperada en tablas interactivas, visualizaciones con gráficos y grafos de coocurrencia de palabras.

La aplicación se diseña como un sistema distribuido bajo la arquitectura cliente-servidor y se soporta en el uso de contenedores, donde en cada uno se ejecuta un proceso para el funcionamiento de la SSCSU, siendo los principales el de la base de datos, el servidor de la aplicación y otro con los distintos procesamientos que son efectuados sobre los textos.

También se propone una rutina que permite clasificar las Tesis y los TEG por el área académica donde cursó estudios el autor del correspondiente trabajo y así se solventa la carencia que actualmente presenta Saber UCV, donde no están disponibles estas clasificaciones.

Índice general

—	1
1 Introducción	2
1.1 Propuesta:	2
1.2 Sistemas de Recuperación de Información:	4
1.3 Saber UCV:	5
1.3.1 Delimitación de documentos con los cuales se trabajará:	5
1.3.1.1 Corpus seleccionado:	6
1.3.1.2 Archivos complementarios:	8
1.4 Requerimiento inicial:	8
1.5 Planteamiento del Problema:	8
1.5.1 I. No se posee clasificación de los documentos:	8
1.5.2 II. Problemas en los resultados que generan los <i>Querys</i> :	9
1.5.2.1 Ejemplo <i>Query</i> en Saber UCV:	9
1.6 Antecedentes:	11
1.7 Estructura del Trabajo:	11
2 Marco teórico-referencial:	13
2.1 Reseña histórica:	13
2.2 Recuperación de Información:	15
2.3 Sistemas de Recuperación de Información (SRI) :	15
2.4 Ejemplos de IRS:	16
2.5 Modelos de Recuperación de Información:	16
2.5.1 Recuperación booleana:	16
2.5.2 Índices Invertidos:	18
2.5.3 Scoring Model:	19

2.6	Procesamientos de texto:	19
2.6.1	Procesamiento del Lenguaje Natural (natural language processing- NLP):	20
2.6.1.1	Tokenizador:	20
2.6.1.2	Entidades Nombradas (<i>named entity reconigition-NER</i>):	20
2.6.1.3	Etiquetado de Partes del Discurso (<i>Part of speech tagging-POS</i>):	20
2.6.1.4	Stemming:	21
2.7	Minería de Texto:	21
2.7.1	Coocurrencia de Palabras:	21
2.8	Similitud de documentos:	22
2.9	Sistemas Distribuidos:	23
2.9.1	Contenedores:	23
2.9.2	Orquestador:	24
2.10	Tendencias actuales Sistemas de Información:	24
2.10.0.1	Temas en proceso de investigación:	25
3	Propuesta Técnica:	26
3.1	Propuesta Técnica:	26
3.2	Objetivo:	27
3.3	Objetivos Específicos	27
3.4	Esquema del SSCSU:	28
3.4.1	Cliente - Servidor:	28
3.4.1.1	1 - Cliente - Navegador web:	28
3.4.1.2	2 - Servidor:	28
3.4.2	Contenedores:	30
3.4.2.1	Docker Compose:	30
3.4.2.2	Contenedor con Nginx:	30
3.4.2.3	Contenedor con Cerbot:	30
3.4.2.4	Contenedor con Shinyproxy:	30
3.4.2.5	Contenedor con “ <i>R Shiny Web App framework</i> ”:	31
3.4.2.6	Contenedor con PostgreSQL:	32
3.4.2.7	Contenedor con “ <i>R Imagen Servicios</i> ”:	33
3.4.2.8	Contenedor con “ <i>Python Spacy</i> ”:	34
3.5	Factibilidad:	34
3.6	Cronograma:	34

Índice de figuras

1.1	Modelo de ficha registro en Saber UCV	6
1.2	Cantidad de documentos disponibles por año de elaboración	7
1.3	Búsqueda de un autor "Jesus Fajardo" en www.saber.ucv.ve	10
1.4	Búsqueda de "simultánea de múltiples metástasis" en www.saber.ucv.ve	10
2.1	Gráfico que acompaña resultados de búsqueda de un término en la biblioteca digital de la Association for Computing Machinery (https://dl.acm.org/)	17
2.2	Coocurrencia de Palabras	22
3.1	Esquema General	28
3.2	Diagrama de Componentes	29

Índice de cuadros

1.1	Total documentos por jerarquía	7
3.1	Cronograma de desarrollo de la Solución	35

Every important aspect of programming arises somewhere in the context of sorting or searching.

— Donald Knuth, *The Art of Computer Programming*, Volume 3

Capítulo 1

Introducción

En este Capítulo se presenta la Propuesta 1.1 denominada **Solución Sistema Complementario Saber UCV (SSCSU)**. Posteriormente se introducen los **Sistemas de Recuperación de Información** 1.2 y el repositorio institucional de documentos electrónicos **Saber UCV** 1.3. Seguidamente se delimitan los documentos de Saber UCV con los que se trabajará en 1.3.1.1. En 1.4 se expone el requerimiento inicial efectuado dentro del Postgrado de Ciencias de la Computación que motivó la realización de este desarrollo y se muestran los impedimentos y limitaciones en 1.5, dadas las características de **Saber UCV**, que era necesario superar para cumplir con el requerimiento y poder desarrollar esta Propuesta. Igualmente se hace un repaso a investigaciones llevadas a cabo 1.6 en la Universidad Central de Venezuela, que parcialmente han abordado problemas y soluciones similares a la que acá se presenta. Finalmente se detalla la estructura 1.7 por capítulos de este Anteproyecto.

1.1 Propuesta:

En este trabajo se formula la propuesta de una Solución que permite realizar procesos de **Recuperación de Información** 2.2 sobre los Resúmenes, tanto de las Tesis como de los Trabajos Especiales de Grado, que están alojados en el repositorio institucional de documentos digitales “Saber U.C.V.” de la Universidad Central de Venezuela.

La propuesta se basa en el desarrollo e implementación de un sistema distribuido con una arquitectura cliente-servidor, disponiendo de una aplicación web donde el cliente accede usando el navegador web. Al usuario se le permitirá hacer la formulación de una consulta para filtrar, dentro de un conjunto de textos, aquellos que cumplen las condiciones definidas.

La consulta, a la cual denominaremos el **query**, estará definida por: una frase, la selección de diversos criterios para el filtrado de los datos y la selección de la granularidad ¹, eg. un intervalo

¹ En la inteligencia de negocios el término granularidad hace referencia al nivel de detalle, o sumariaización de las unidades de datos que reposan en la *data warehouse*, a más detalle, menos es la granularidad y a menor detalle implica una mayor granularidad (Inmon, 2002). En el campo de la recuperación de información este término se contextualiza en el nivel de detalle al que llegará la recuperación dentro de un texto, pudiendo ser por ejemplo dentro de un documento la búsqueda de una ocurrencia, dentro de un capítulo, de una sección, de un párrafo o de una oración. Una de las formas en que serán representados los resultados obtenidos, mediante la coocurrencia de palabras, mostrará distintas granularidades tanto en párrafos como en el texto completo del resumen de la tesis. Igualmente se asocia la granularidad a la restricción que se podrá seleccionar para hacer la búsqueda, delimitando si se quiere buscar dentro de todas las facultades, en alguna de ellas, o sólo documentos dentro producidos dentro de una escuela y/o postgrado.

de fechas, la selección de la(s) facultad(es) y/o escuela(s)-postgrado(s) donde se generó la investigación. Al ejecutar la consulta serán recuperados del Sistema los textos que cumplen con las restricciones establecidas. Los resultados se mostrarán con distintos gráficos, tablas interactivas, grafos de coocurrencia de palabras 2.7.1 y recomendaciones de documentos, basados en las similitudes que pueda presentar un documento con los otros 2.8. La propuesta descrita recibirá el nombre de **Solución Sistema Complementario Saber UCV (SSCSU)**.

La SSCSU será desarrollada basándose en los principios que definen a los Sistemas de Recuperación de Información 2.3 y se conformará un conjunto delimitado de documentos con los resúmenes, lo cual en la literatura especializada se denomina el *Corpus* 1.3.1.1. Sobre este conjunto serán aplicadas técnicas de Minería de Texto 2.7 y de Procesamiento de Lenguaje Natural 2.6.1, para detectar relaciones entre los elementos y extraer conocimiento de los mismos (Aggarwal and Zhai, 2012a).

La implementación del Sistema se hará mediante el uso de contenedores encargados de sustentar los procesos modulares y complementarios. Entre los procesos modulares se encuentran los de gestión de la base de datos, el servidor de la aplicación y el de procesamientos de los textos ². En la propuesta técnica 3.1 se darán los detalles de las procesos que tendrán encargados ejecutar cada uno de los contenedores.

Uno de los aspectos relevantes a resolver para implementar este Sistema, es realizar la clasificación de las Tesis por área de conocimiento, asociándose a la facultad, escuela o el postgrado, donde cursó estudios el autor de cada investigación. Se destaca que actualmente esta información no está disponible para el público ni en el repositorio Saber UCV. Se pretende que el proceso de clasificación sea realizado con una elevada precisión, derivando en un beneficio para la comunidad de investigadores.

Previo a la ejecución de lo antes descrito se realizará una investigación documental para generar un listado con todas las etiquetas candidatas con las que se ejecutará el proceso de clasificación. Para esto se usarán principalmente los listados de carreras de pre y post grados disponibles en las páginas oficiales de la Universidad y sus facultades. La extracción de las etiquetas se hará valiéndose de identificadores en el código HTML y CSS de estas. El procedimiento propuesto implica la facilidad de poder reproducir el código de extracción de tales valores y hacer futuras actualizaciones, siempre y cuando se mantenga la estructura en las páginas web mencionadas.

El proceso de clasificación, adicional a la necesidad de contar con las etiquetas, lleva a que se tengan que descargar los archivos anexos a cada uno de los Resúmenes, estando estos documentos en formato *Portable Document Format* (PDF), word, entre otros. Una vez descargados se hace la lectura de estos archivos, resolviendo otra diversidad de problemas como las codificaciones (*file encodings*) de estos, el tipo de información que pueden tener (texto o imagen) y que algunos resúmenes están compuestos por más de un archivo, haciendo necesario la consolidación de los mismos. Todo esto finalmente lleva a la extracción de un trozo de texto para poder lograr el emparejamiento con la etiqueta correcta. Lo antes descrito lleva a que sea necesario contar con diversas herramientas de procesamientos de archivos y algoritmos de tratamiento de textos, para lograr el objetivo de ejecutar la clasificación con la debida precisión.

También la SSCSU permitirá obtener indicadores cuantitativos sobre los trabajos que reposan en Saber UCV, pudiendo establecer distintos criterios de filtrado, para la obtención de estas cifras. Entre estos indicadores podemos mencionar la cantidad de documentos que cumplen con el

²La creación de los contenedores y la integración de estos se hará mediante el uso de *Docker* y *Docker Compose* respectivamente

1.2. SISTEMAS DE RECUPERACIÓN DE INFORMACIÓN:

criterio del *query* por año de elaboración de la Tesis o el TEG, o la cantidad de documentos totales recuperados por área del conocimiento: facultad, postgrado y escuela.

La SSCSU contará con una base de datos estructurada indexada mediante un “tsvector” el cual es un tipo de datos nativo del manejador de base de datos que usaremos (PostgreSQL) y que permite representar y guardar los documentos de una forma optimizada para la búsqueda de texto. En 2.5.2 se desarrolla el genérico de este tipo de estructura.

El Sistema propuesto tiene que poder acoplarse periódicamente a Saber U.C.V. para realizar la descarga, clasificación e incorporación al *Corpus* de los nuevos documentos que estén disponibles en el repositorio oficial y por ello estará diseñado para actualizarse periódicamente, incorporando así los nuevos trabajos.

Un aspecto que también se toma en cuenta en el diseño de la SSCSU es la usabilidad y la experiencia del usuario al interactuar con él. Para esto, previendo la complejidad que pueda representar el uso del Sistema, en la interfaz gráfica que visualiza el usuario, se incorporarán ventanas emergentes que sirvan para explicar las funcionalidades de cada una de las áreas de interacción.

A modo resumen, el objetivo de esta investigación es construir un sistema de recuperación de información que funcione bajo una arquitectura distribuida cliente-servidor, que permita clasificar los trabajos que reposan en Saber UCV y a los usuarios finales generar consultas de búsqueda de texto sobre los documentos que constituyen el *Corpus* mediante el acceso a una aplicación web.

1.2 Sistemas de Recuperación de Información:

Es conocida la gran diversidad de información en distintos formatos que tienen a su disposición los investigadores. Ejemplo de esto son libros en las bibliotecas o los documentos que se encuentran accesibles en formatos digitales como artículos publicados en revistas arbitradas, libros, páginas de internet especializadas en algún tema o los repositorios digitales. Es en esta abundancia donde recaen varios de los problemas a los cuales se enfrenta la labor de investigación (Hernández Orallo et al., 2004) (p.565), especialmente cuando se realiza la fase exploratoria de selección de aquellos documentos que resultan de mayor interés.

El investigador debe contar con herramientas al enfrentarse al proceso de búsqueda de información siendo una de estas los ***Sistemas de Recuperación de Información (SRI)***, que son los que permiten que el investigador formule, con el nivel de detalle deseado, la búsqueda mediante un texto al cual denominaremos *query*, siendo el sistema el encargado de detectar la aparición de dicho texto en los documentos que previamente han sido agrupados como un *Corpus* (Manning et al., 2008).

Conocemos que tal búsqueda se debe ejecutar con ciertos niveles de flexibilidad y no restringirse a la aparición exacta del *query*, permitiendo que en vez de localizar la aparición exacta de una, o varias palabras, sea la raíz de las mismas la que se localice, eg. investigamos, investigan, investigaron, tendrá como raíz “*investig*”. Caso similar se puede presentar con una palabra que se indique en el *query* en plural pero al momento de efectuar la búsqueda, sea el singular el que también sea detectado. Los dos puntos anteriores lo que hacen es ensanchar el espacio de búsqueda para así incrementar las posibles fuentes de información que serán recuperadas y presentadas al investigador.

Por otra parte el espacio de búsqueda se puede estrechar filtrando los documentos que estén comprendidos entre un período restringido de fechas o que esté asociado a una determinada clasificación, como por ejemplo, a un área del conocimiento específica.

Para ilustrar lo expuesto en los precedentes párrafos usemos el siguiente ejemplo: supongamos que un ingeniero quiere buscar información sobre la “investigación de operaciones” dentro de un **Corpus** de textos académicos que contiene diversas áreas de conocimiento. Al hacer el *query* el sistema que usemos va a determinar cuáles documentos dentro del *corpus* hacen mención a ese tema. Es probable que encontrásemos textos asociados a medicina donde la “investigación de operaciones” tiene una semántica distinta a la rama del conocimiento que estudia problemas de optimización que se denomina “investigación de operaciones”. Esto es gracias a que en diferentes áreas de estudio se pueden tener comprensiones que divergen sobre una misma búsqueda de términos.

Visto lo anterior, es necesario que el Sistema de Recuperación de Información que usemos permita separar de alguna forma la información recuperada, acorde a un contexto de interés del investigador, jerarquizando aquellos documentos que puedan resultar más provechosos. Incluso, las visualizaciones con las que representemos los resultados obtenidos de un *query*, pueden aportar conocimiento en sí mismo haciendo intuitiva la comprensión (Zhang, 2008) o dando aportes hacia dónde es viable dirigir futuras indagaciones, en un proceso que claramente es iterativo (Zhai and Massung, 2016).

Adicionalmente el sistema puede generar la recomendación de otros trabajos (Aggarwal, 2018) que presenten alguna similitud sobre el tema investigado. En la sección 2.3 se exponen otros detalles relevantes sobre los **SRI**.

1.3 Saber UCV:

La Universidad Central de Venezuela cuenta con un Sistema de Recuperación de Información denominado **Saber UCV**, al cual se puede acceder en la dirección www.saber.ucv.ve que funciona a modo de repositorio digital institucional ³ de los distintos documentos de textos académicos generados por la comunidad ucevista. Este sistema fue “creado para alojar, gestionar y difundir de manera gratuita y en texto completo: tesis, artículos de investigación, libros, revistas electrónicas, presentaciones que conforman la producción académica” de esta casa de estudios.

Saber UCV es una implementación del software de código abierto llamado **DSpace** que fue creado con la finalidad de ser un repositorio digital institucional bibliográfico, para uso académico de colecciones digitales según se indica en la página web oficial del proyecto <https://duraspace.org>.

Este software fue diseñado e implementado en un esfuerzo conjunto entre desarrolladores del *Massachusetts Institute of Technology (MIT)* y *Hewlett-Packard Labs (HP Labs)* en el año 2002. En junio de 2021 se lanzó la 7 versión que es la más reciente, no obstante la versión que está implementada en Saber UCV al mes de noviembre del 2021 es la 1.7.1. que data del año 2013.

³“base de datos electrónica que aloja una colección de pequeñas unidades de información educativas o actividades que pueden ser accedidas para su obtención y uso” (Lehman, 2007)

1.3.1 Delimitación de documentos con los cuales se trabajará:

Si bien en el sitio web están almacenados artículos de prensa, artículos preimpresos y publicados, fotografías, infografías y tesis de distintos niveles académicos, es exclusivamente con este último subconjunto con el cual se trabajará.

1.3.1.1 Corpus seleccionado:

En la categoría **Tesis** del repositorio se encuentran los trabajos especiales de grado de pre y postgrado junto con las tesis de doctorado. Cada documento de la categoría Tesis, indistintamente del nivel académico al que pertenezca, tiene uno o más archivos anexos en formato *word* o *PDF* donde está contenida la investigación y adicionalmente cada documento cuenta con una ficha en *HTML*. En esta investigación se usará de ahora en adelante la denominación “Tesis”, indistintamente de si son trabajos especiales de grado de pre y postgrado o tesis de doctorado.

En esta ficha se muestran ciertos detalles sobre el trabajo seleccionado como lo son: el título, el nombre del autor, las palabras claves, la fecha de publicación y el texto **Resumen**. En la figura 1.1 podemos ver la ficha que corresponde a un trabajo almacenado en el repositorio.

Con los textos **Resumen**, con la **fecha**, el **nombre del autor**, las **palabras claves** de todos los trabajos catalogados como Tesis en el repositorio Saber UCV es con los que se conformará el **Corpus**.

Para el momento en que se presenta esta Investigación bajo la categoría **Tesis** reposan 10.481 documentos distribuidos con las jerarquías que se muestran en el cuadro 1.1.

Cuadro 1.1: Total documentos por jerarquía

jerarquía	cantidad
doctorado	305
maestría	729
otras	1.239
pregrado	8.208
Total	10.481

La disponibilidad de trabajos en el repositorio por fecha de creación la presentamos en el histograma que se muestra en la figura 1.2 .

Analizando la distribución de la frecuencia de la figura 1.2 por año en que fue generada la investigación es necesario destacar que en el repositorio no se encuentran todas las investigaciones que se han realizado en la historia de la U.C.V. Actualmente se tienen incorporados solo 20 trabajos creados entre 1993 y 2011, siendo a partir de este último año cuando se empieza a realizar la ingesta de forma periódica y en mayor volumen. No obstante también queremos resaltar que desde el 2011 no se encuentran incorporados todos los trabajos generados dentro de la comunidad. Las

Por favor, use este identificador para citar o enlazar este ítem:
<http://hdl.handle.net/10872/12302>

Título :	Desarrollo de un Método de Ordenamiento Vectorial para la extensión de los Operadores Básicos en Morfología Matemática a imágenes Multidimensionales
Autor :	Acevedo P, Yorley A,
Palabras clave :	operaciones morfológicas imágenes escala de gris imágenes binarias imágenes multidimensionales ordenamiento vectorial
Fecha de publicación :	16-Oct-2015
Resumen :	El presente Trabajo Especial de Grado se orienta al estudio e investigación de la extensión de las operaciones morfológicas básicas que se aplican a imágenes en escala de gris e imágenes binarias a imágenes multidimensionales o vectoriales que poseen un mismo tipo de información en cada uno de sus canales. El principal aporte de este trabajo es proponer un nuevo método de ordenamiento vectorial para la aplicación de operaciones morfológicas en imágenes vectoriales o multicanales (imágenes a color, resonancia magnética, rayos X, satelital, etcétera).
URI :	http://hdl.handle.net/10872/12302
Aparece en las colecciones:	Pregrado

Ficheros en este ítem:

Fichero	Descripción	Tamaño	Formato	
Tesis Yorley A. Acevedo P..pdf		8.71 MB	Adobe PDF	Visualizar/Abrir

Figura 1.1: Modelo de ficha registro en Saber UCV

1.4. REQUERIMIENTO INICIAL:

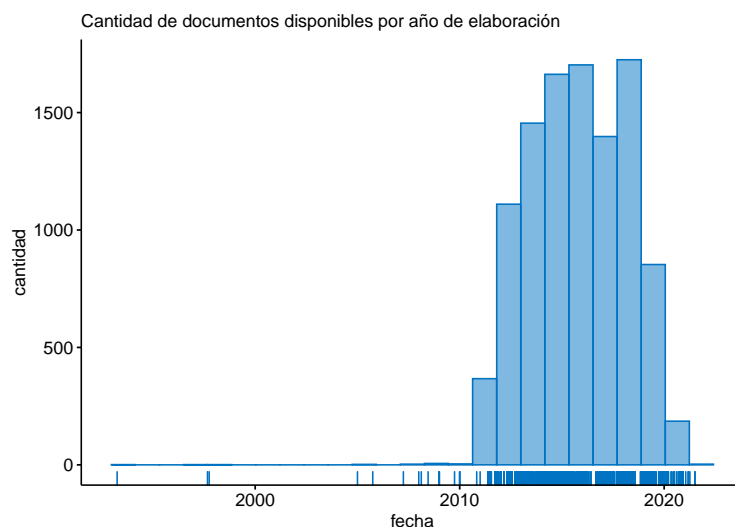


Figura 1.2: Cantidad de documentos disponibles por año de elaboración

razones que explican las omisiones o retardos en las incorporaciones queda fuera de los objetivos de este trabajo investigarlas.

1.3.1.2 Archivos complementarios:

Por razones que se expondrán en 1.5.1 igualmente será necesario trabajar con los documentos anexos que indicamos en 1.3.1.1 que disponen un vínculo de descarga en cada ficha *HTML*. La descarga y el procesamiento de estos textos será necesario para complementar el corpus descrito en 1.3.1.1.

1.4 Requerimiento inicial:

Originalmente en la materia *Tópicos Especiales en Sistemas de Información y Gerencia* se hizo un requerimiento para con los trabajos especiales de grado realizados por los alumnos de la Facultad de Ciencias, disponibles en Saber UCV, se hicieran *Mapas de Conocimiento*.

Estos Mapas en un trabajo de (Dueñas et al., 2011) se definen como “una metodología que busca encontrar palabras claves como unidad de análisis que permitan establecer cómo se relacionan diferentes campos científicos alrededor de un tema determinado”.

Ante esto se buscaron las estrategias que pudiesen satisfacer el requerimiento. Lo primero que se necesitaba era conformar el conjunto de datos, ya que no fue entregado como insumo.

Al realizar la descarga de los documentos con un *web crawler* de diseño propio se pudo apreciar que no era viable separar los documentos de la Facultad de Ciencias al no existir ninguna etiqueta para realizar la selección.

Esto derivó en que fuese necesario diseñar una solución para clasificar los documentos previo a poder generar los Mapas de Conocimiento. En 1.5.1 se expone con mayor detalle este problema.

1.5 Planteamiento del Problema:

El repositorio SABER UCV permite realizar búsqueda de documentos mediante procesos de recuperación de información, no obstante hay algunas características que pueden ser complementadas y mejoradas. A continuación se listan.

1.5.1 I. No se posee clasificación de los documentos:

El **Corpus** que reposa en Saber UCV para ninguno de los documentos posee la etiqueta que clasifique a cual postgrado o carrera de pregrado pertenece. En caso de realizar una búsqueda y querer delimitarla a documentos que se circunscriban a una determinada área de conocimiento, actualmente es imposible y por ende resulta inviable evaluar la escuela o el postgrado y la facultad donde se hizo la investigación.

Tampoco se dispone de un listado con las etiquetas candidatas a asignar a cada uno de los documentos. Según lo investigado no es de conocimiento público un único directorio que contenga los nombres de las carreras de pregrado y de los distintos postgrados, significando esto que es un doble proceso de construcción el que se necesita hacer para realizar la clasificación de los documentos.

Relativo a este problema también surge el determinar cómo se realizará la comparación de cada una de las etiquetas con el texto de cada documento, incluso en detectar dónde se localiza el texto con el cual se efectuará la comparación, ya que el texto denominado **Resumen**, que vimos en la figura 1.1, no contiene la información necesaria para extraer dónde fue realizada la investigación (escuela o postgrado), haciendo necesario que sean analizados los documentos anexos a cada ficha que es propiamente el archivo, o los archivos, que componen cada una de las investigaciones.

En detalle, se tendrán que consolidar los archivos disponibles para cada investigación, los cuales no necesariamente se encuentran ordenados, por ejemplo por capítulo. Una vez realizada la consolidación se procederá a crear un archivo en texto plano. En algunos casos los documentos no contendrán texto sino imágenes del texto motivando que se tenga que realizar lectura OCR ⁴ para extraer el texto.

En este proceso de clasificación también se presentarán otros problemas relativos a los sistemas de codificaciones de los archivos *file encodings* y a la introducción de caracteres no deseados dentro del texto, lo que será expuesto a detalle en el capítulo que contiene el desarrollo de la propuesta.

1.5.2 II. Problemas en los resultados que generan los *Queries*:

Existe un tipo de *query* denominado “*Full Text Search*” o búsqueda de texto completo que es donde se manifiesta una de las principales deficiencias que posee el Sistema Saber UCV. Aunque en él se pueden ejecutar búsquedas de texto, los métodos con los que se implementan estas búsquedas llevan a que sean generados una gran cantidad de resultados, en especial al usar más de una palabra en la búsqueda, presumimos que por usar el operador lógico *OR* entre cada una de las palabras solicitadas en el *query*.

⁴OCR: “*optical character recognition*” traducido como reconocimiento óptico de caracteres

1.5. PLANTEAMIENTO DEL PROBLEMA:

1.5.2.1 Ejemplo *Query* en Saber UCV:

Evaluemos el problema descrito con un ejemplo. Supongamos que queremos ver en el sistema Saber UCV los trabajos que un autor de nombre “Jesús Fajardo” ha realizado.

The screenshot shows the 'Resultados de búsqueda' (Search Results) page on the Saber UCV website. The search criteria are set to 'Buscar: 2) Tesis' and 'por: |jesus fajardo|'. The results show 1-10 of 30 items. The first three results are displayed in a table:

Fecha de publicación	Título	Autor(es)
13-Nov-2017	Marca personal para estudiantes de comunicación social : proyecto audiovisual sobre la importancia de construir y gestionar una marca personal.	Bulló Camejo, Vanessa Valentina; Fajardo García, Karen Anamileth
24-Nov-2015	Irradiación simultánea de múltiples metástasis cerebrales con técnica de vmat. Comparación con técnica de radiocirugía	Fajardo Freitas, Jesús Ernesto
27-Nov-2019	Capacidad de carga para la gestión del sector Sabas Nieves del parque nacional Waraira Repano	Fajardo, Eilezer

Figura 1.3: Búsqueda de un autor “Jesus Fajardo” en www.saber.ucv.ve

Al revisar los resultados que se muestran en la figura 1.3 podemos apreciar que la primera posición de los 30 resultados que fueron encontrados, no muestra un autor que se llame “Jesús Fajardo” sino el apellido Fajardo. Incluso, al realizar la consulta del vínculo del primer resultado tampoco aparece la palabra “Jesús” dentro de todos los datos que muestra la ficha. Al revisar el segundo resultado vemos que sí coincide con el término que se estaba buscando, pero de resto ningún otro de los 28 resultados incluyen un “Jesús Fajardo” o un “Jesús”.

Revisemos otro ejemplo de búsqueda esta vez con cuatro palabras que son “simultánea de múltiples metástasis” las cuales se encuentran en el título del trabajo visto anteriormente en la posición número dos. El resultado del *query* lo podemos ver en la figura 1.4.

La cantidad de resultados asciende a 10.449, cifra similar a la cantidad total de trabajos de la categoría “Tesis” que se encuentran en el repositorio, ya que como fue explicado anteriormente el Sistema Saber UCV debe aplicar el operador lógico *OR* entre cada una de las palabras al momento de generar el *query*. A manera de referencia podemos explicar que si hubiésemos querido realizar la búsqueda exacta de la frase, con haberle colocado comillas al inicio y al final obtendríamos la referencia al trabajo que previamente sabemos que es el que queremos encontrar.

Sin embargo, el uso de comillas para buscar una frase exacta genera una compensación entre la precisión en la búsqueda, lo cual se motiva en que al ser exacta, se cierra el umbral para encontrar otros resultados que pudiesen resultar de interés, así que el uso de este método puede presentar considerables desventajas.

Posterior a evaluaciones realizadas sobre el comportamiento de Saber UCV presumimos que los textos al ser procesados para su estructuración e ingesta en la base de datos no le son aplicadas

Buscar: 2) Tesis

por simultánea de múltiples Ir

Resultados 1-10 de 10449.

Resultados por página 10 | Ordenar por Relevancia En orden Descendente Autor/registro Todo Actualizar

Resultados por ítem:

Fecha de publicación	Título	Autor(es)
24-Nov-2015	Irradiación simultánea de múltiples metástasis cerebrales con técnica de vmat. Comparación con técnica de radiocirugía	Fajardo Freites, Jesús Ernesto

Figura 1.4: Búsqueda de "simultánea de múltiples metástasis" en www.saber.ucv.ve

técnicas con el algoritmo de Porter que permite realizar la extracción de la raíz o (*stemming*)⁵ de cada una de las palabras. La importancia de aplicar este procesamiento será explicado en 2.6.1.4.

1.6 Antecedentes:

Se hizo una revisión bibliográfica de trabajos producidos dentro de la Universidad Central de Venezuela que puedan haber realizado investigaciones similares o complementarias a la que proponemos en este desarrollo. De particular interés resultó la "Elaboración de un prototipo de buscador de documentos académicos de la Facultad de Ciencias" (Sánchez, 2008), que fue una investigación que implementó el motor de búsqueda y repositorio denominado **BUSCONEST 1** que sirve exclusivamente a los documentos de investigación generados en la Facultad de Ciencias de la U.C.V. Este sistema sí cuenta con los trabajos clasificados por área del conocimiento donde fueron generados, no por contar con una rutina para la clasificación o por disponer de una etiqueta, sino por ser otra la fuente de los datos, en este caso el sistema de datos de la Biblioteca de la Facultad de Ciencias. Con lo analizado se comprende que no se resuelve el problema de clasificar el resto de documentos que reposan en Saber UCV.

Posteriormente (Guevara, 2015) presenta otra propuesta de motor de búsqueda con una versión modificada de BUSCONEST 1, denominada BUSCONEST 2, no obstante sólo se basa sobre textos producidos en la Facultad de Ciencias y en ninguno de estos desarrollos se realizan procesos de minería de texto.

Otra investigación que podemos señalar es "Sistema de recomendación para el Buscador Académico Venezolano" (Rodríguez Laguna, 2016). En ella se señalan algunos repositorios digitales institucionales que se encuentran en Venezuela y se hace una propuesta de unificación de todos los contenidos y así poder generar recomendaciones de investigaciones a partir de una determinada búsqueda.

⁵proceso heurístico que remueve letras al final de las palabras con el objetivo de encontrar la raíz de la misma.

1.7. ESTRUCTURA DEL TRABAJO:

Evalando estos antecedentes consideramos que nuestra propuesta resuelve el problema de la clasificación de documentos y permite efectuar las búsquedas sobre todo el conjunto de textos, así que podemos tener en perspectiva que nuestro desarrollo no solapará las propuestas mencionadas y generará un valor agregado favoreciendo a la actividad de investigación.

1.7 Estructura del Trabajo:

Este trabajo está estructurado de la siguiente forma: en un Capítulo 2 se expone el Marco Teórico. En el Capítulo 3 se describirá la propuesta técnica, el objetivo general y específicos de la Solución.

Tentativamente se propone que en el Proyecto de Trabajo de Grado en el Capítulo 4 se exponga el proceso de desarrollo de la Solución y las distintas pruebas que se realizaron para elegir una determinada tecnología por sobre otras, así como los resultados obtenidos, mientras que en el Capítulo 5 se presentarán las conclusiones y recomendaciones, posterior a la implementación del Sistema.

Capítulo 2

Marco teórico-referencial:

En este Capítulo se muestra el marco teórico en que se sustentan los aspectos de mayor relevancia para el desarrollo de la Solución. Principalmente se enuncian una serie de conceptos que involucran algoritmos de búsqueda, la recuperación de información, la minería de texto, el procesamiento del lenguaje natural, la estructuración de la base de datos y lo referente a la arquitectura distribuida en que se soporta la SSCSU.

2.1 Reseña histórica:

El profesor Donald Knuth señala, dentro del campo de las ciencias de la computación, que la **búsqueda** *es el proceso de recolectar información que se encuentra en la memoria del computador de la forma más rápida posible, esto cuando tenemos una cantidad N de registros y nuestro problema es encontrar el registro apropiado de acuerdo a un criterio de búsqueda* (Knuth, 1997) (p. 392) .

Iniciamos con esta cita porque la recuperación de información gira en torno a un problema central de las ciencias de la computación que es la **búsqueda**. A continuación mencionaremos una serie de algoritmos que abordan este problema, no necesariamente resultando óptimos para dar solución a lo planteado en 1.5.2.

En la década de 1940 cuando aparecieron las computadoras, las búsquedas no representaban mayor problema debido a que estas máquinas disponían de poca memoria *RAM* pudiendo almacenar sólo moderadas cantidades de datos. Ellas estaban diseñadas para realizar cálculos y arrojar los resultados más no para tenerlos almacenados en memoria.

No obstante con el desarrollo del almacenamiento en memoria *RAM* o en dispositivos de almacenamiento permanentemente, ya en la década de 1950 empezaron a aparecer los problemas de **búsqueda** y los primeras investigaciones para afrontarla. En la década de 1960 se adoptan por ejemplo estrategias basadas en árboles.

Los primeros algoritmos que sirvieron para localizar la aparición de una frase dentro de un texto, o expresado de forma más abstracta, como la detección de una subcadena P dentro de otra cadena T , fueron los algoritmos de *Pattern-Matching* (Goodrich et al., 2013) (p. 584).

2.1. RESEÑA HISTÓRICA:

Así nos encontramos en la literatura con el algoritmo *Fuerza Bruta* donde dado un texto T y una subcadena P, se va recorriendo cada elemento de la cadena T para detectar la aparición de la subcadena P. Si bien este algoritmo no presentaba el mejor desempeño, por contar con ciclos anidados en su ejecución, creó una forma válida de enfrentar el problema de la búsqueda de subcadenas de texto.

El algoritmo *Knuth-Morris-Pratt* que se introdujo en 1976 tenía como novedad que se agregó una función que iba almacenando “previas coincidencias parciales” en lo que eran fallos previos y así al realizar un desplazamiento tomaba en cuenta cuántos caracteres se podían reusar. De esta forma se logró considerablemente mejorar el rendimiento en los tiempos de ejecución de $O(n+m)$ que son asintóticamente óptimos.

Posteriormente en 1977 el problema se enfrenta con un nuevo algoritmo que es el de *Boyer-Moore* en el cual se implementan dos heurísticas (*looking-glass* y *character-Jump*) que permiten ir realizando algunos saltos en la búsqueda, ante la no coincidencia de la subcadena con la cadena y adicionalmente, el orden en el que se va realizando la comparación se invierte. Estas modificaciones permitieron obtener un mejor desempeño.

Sobre una modificación al algoritmo *Boyer-Moore* se sustenta la utilidad *grep* de la línea de comandos UNIX que igualmente le da soporte a diversos lenguajes que la usan para ejecutar búsquedas de texto con un proceso que comúnmente es conocido como *grepping*. Esta utilidad fue ampliamente usada para resolver parcialmente lo expuesto en 1.5.1.

Los algoritmos mencionados anteriormente pueden ser usados en procesos de recuperación de información en conjunto con técnicas que pueden mejorar considerablemente los tiempos en la ejecución de las rutinas siendo una de estas el preprocesamiento de los textos, *eg.* remover determinados caracteres, aplicar el algoritmo de porter, entre otras más.

Una estrategia que surgió para enfrentar las búsquedas de texto, fue el uso de la programación lineal donde bajo la premisa *divida et impera* los problemas que requieren tiempo exponencial para ser resueltos son descompuestos en polinomios y por lo tanto se disminuye la complejidad en tiempo para ser resueltos. Entre este tipo de algoritmos podemos mencionar los de *alineación del ADN*. Originalmente el algoritmo se desarrolló para resolver problemas de alineación de cadenas de ADN de forma parcial o total dentro de una cadena mayor. Posteriormente se identificó que este tipo de procedimiento era extrapolable a los subcadenas de texto. Una de las versiones de estos algoritmos es la denominada ***Smith-Waterman*** que resultó de gran utilidad para resolver el problema planteado en 1.5.1 ya que la estrategia de usar la utilidad *grep* fue infructuosa en algunos casos.

Un gran paso para aproximarnos a la aparición de los Sistemas de Recuperación de Información lo representó el enfoque que presentan los algoritmos *Tries*. Este nombre proviene del proceso de *Information Retrieval* y principalmente se basa en hacer una serie de preprocesamientos a los textos para que al momento de ejecutar la búsqueda de texto, es decir, de la subcadena dentro de la cadena, ya tengamos una parte del trabajo realizado previamente y no tener que ejecutarlo todo “*on the fly*”, es decir, sobre la marcha.

Un *Trie* (Fredkin, 1960) es una estructura de datos que se crea para almacenar textos para así poder ejecutar más rápido la coincidencia en la búsqueda. En la propuesta de la SSCSU todos los textos van siendo procesados con distintas técnicas a medida que son insertados en la base de datos. Esto claramente representa una mejora en el desempeño con la disminución en los tiempos de búsqueda.

2.2 Recuperación de Información:

El eje central sobre el cual gira el proceso de recuperación de información (RI) es satisfacer las necesidades de información relevante que sean expresadas por un usuario mediante una consulta (*query*) de texto. El investigador Charu Aggarwal en su libro sobre Minería de Texto (Aggarwal and Zhai, 2012b) (p.14) menciona que el objetivo del proceso de RI es conectar la información correcta, con los usuarios correctos en el momento correcto, mientras que otro de los autores con mayor dominio sobre el tema, Christopher Manning en su libro *Information Retrieval* indica que “es el proceso de encontrar materiales (generalmente documentos) de una naturaleza no estructurada (generalmente texto) que satisface una necesidad de información dentro de grandes colecciones (normalmente almacenada en computadores)” (Manning et al., 2008) (p. 25).

Satisfacer una necesidad de recuperación de información no sólo se circunscribe a un problema **búsqueda** de un texto dentro de un *corpus*. En la mayoría de los casos se deberá cumplir con ciertos criterios, o restricciones, como por ejemplo que el *query* esté dentro de un período de fechas, o que se encuentre comprendido en un subconjunto del corpus que es a lo que se denomina **búsqueda múltiple atributo**.

La información que se recolecte en una búsqueda tendrá distintos aspectos que aportarán peso en el orden en que sea presentada al usuario y no sólo vendrá dado por la aparición de las palabras sino también por otros elementos como lo puede ser la aparición de la frase del *query* dentro del título, la proximidad (la cercanía entre dos palabras) que tengan los términos que conforman el *query* o por otra parte la frecuencia que una palabra, o varias, se repitan dentro un determinado documento que compone el *corpus*. Igual puede aportar un peso mayor a la recuperación de un documento las referencias (citas) que contengan otros documentos a ese determinado documento. El fin último será la extracción de los documentos que resulten de mayor relevancia para el usuario.

Incluso es válido incorporar documentos, en los resultados que arroje la búsqueda, que propiamente no coincidan con los términos buscados sino que contengan términos que sean sinónimos o también añadiendo a los resultados documentos que presenten alguna similitud con el texto del *query*. Lo que acabamos de mencionar incorporará formalmente dentro del proceso de extracción de información algo de imprecisión con la intención última de enriquecer el proceso de *Information Retrieval* (Kraft and Colvin, 2017).

Evaluando el proceso con cierto nivel de abstracción tenemos que el proceso de recuperación de información está compuesto principalmente por: un *query*, por un corpus y por una función de *ranking* para ordenar los documentos recuperados de mayor importancia a menor.

El desarrollo de los algoritmos expuestos en 2.1 sumado a la necesidad de resolver los problemas asociados a la búsqueda de un texto dentro de un *corpus* con múltiples atributos en tiempos aceptables y la abundante cantidad de información digital, potenciada por el uso generalizado de los computadores, abonó las condiciones para la creación de los **Sistemas de Recuperación de Información**.

2.3 Sistemas de Recuperación de Información (SRI) :

Los Sistemas de Recuperación de Información (*Information Retrieval Systems-IRS*) son los dispositivo (software y/o hardware) que median entre un potencial usuario que requiere información

2.4. EJEMPLOS DE IRS:

y la colección de documentos que puede contener la información solicitada (Kraft and Colvin, 2017) 1. El SRI se encargará de la representación, el almacenamiento y el acceso a los datos que estén estructurados y se tendrá presente que las búsquedas que sobre él recaigan tendrán distintos costos siendo uno de estos el tiempo que tarde en efectuarse.

Es de nuestro conocimiento que generalmente los datos estructurados son gestionados mediante un sistema de base de datos pero en el caso de los textos estos se gestionan por medio de un motor de búsqueda motivado a que los textos en un estado crudo carecen de estructura (Aggarwal and Zhai, 2012b) (p. 2). Son los motores de búsqueda (*search engines*) los que permiten que un usuario pueda encontrar fácilmente la información que resulte de utilidad mediante un *query*.

El SSCSU está diseñado como un IRS donde se pueden ejecutar queries que son procesados y los resultados que se obtienen son sometidos a una función de ranking que será expuesta en una fase posterior del desarrollo de esta investigación.

2.4 Ejemplos de IRS:

Profundizando en el tema de esta Investigación mencionaremos un par de páginas de internet que funcionan como IRS.

1. Una es el proyecto denominado Arxiv alojado en <https://arxiv.org/>, que es un repositorio de trabajos de investigación. Al momento del usuario hacer un requerimiento de información, adicional al texto de la búsqueda, se pueden indicar distintos filtros a aplicar como puede ser el área del conocimiento (física, matemática, computación, etc.). Igualmente se puede indicar si se quiere buscar sólo dentro del título de una investigación, o el autor, en el *abstract*, o en las referencias, por ejemplo.

Al ejecutar una búsqueda pueden ser recuperados miles de documentos y la interacción con el sistema permite ver que se genera un *ranking* en la exhibición de los resultados obtenidos. El primero de estos *rankings* se ordena con base en la fecha de publicación, pero es viable que se ordenaran los documentos por la relevancia que presentan.

2. También tenemos el portal de la Association Computery Machine (ACM) que incorpora motores de búsqueda con particulares características facilitando la labor de investigación y extracción de información ante una determinada necesidad. Esto lo decimos porque los resultados de una búsqueda son acompañados de distintas representaciones gráficas que le dan un valor adicional a la representación de los resultados. En la figura 2.1 vemos una de estas representaciones que incluye la frecuencia de aparición del *query* en el tiempo.

2.5 Modelos de Recuperación de Información:

2.5.1 Recuperación booleana:

Ante una búsqueda de información se recorre linealmente todo el documento para retornar un valor booleano indicando la presencia o no del término buscado. Es uno de los primeros modelos que se

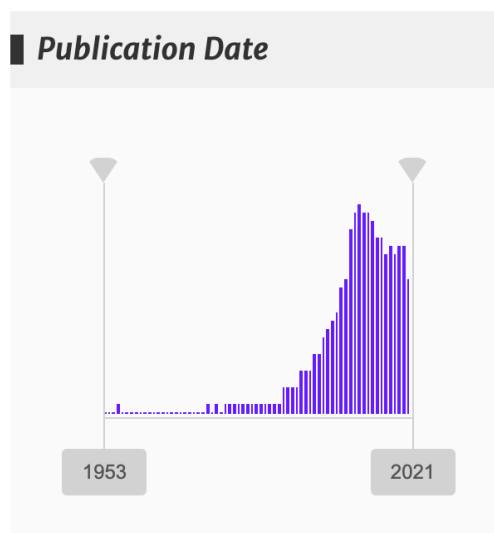


Figura 2.1: Gráfico que acompaña resultados de búsqueda de un término en la biblioteca digital de la Association for Computing Machinery (<https://dl.acm.org/>)

uso y está asociado a técnicas de *grepping* (Manning et al., 2008) (p.3). El desarrollo de este modelo apareció entre 1960 y 1970.

El usuario final obtendrá como respuesta a su *query* sólo aquellos textos que contengan el término. Es un modelo muy cercano a los típicos *queries* de bases de datos con el uso de operadores “AND”, “OR” y “NOT”. En el procesamiento de los textos se genera una matriz de incidencia binaria término-documento, donde cada término que conforma el vocabulario, ocupa una fila i de la matriz mientras que cada columna j se asocia a un documento. La presencia de el término i en el documento j se denotará con un valor verdadero o un “1”.

La recuperación booleana si bien representa una buena aproximación a la generación de *queries* más rápidos, presenta una gran desventaja y es que al crecer la cantidad de documentos y el vocabulario, se obtiene una matriz dispersa de una alta dimensionalidad que hace poco efectiva su implementación.

Los documentos y los *queries* son vistos como conjuntos de términos indexados, que luego fueron llamados “bolsa de términos” (*bag of terms*). Las deficiencias de este modelo recaen en que los resultados, no tienen ningún ranking. Si por ejemplo el término sobre el cual se realiza el *query* aparece 100 veces en un documento y en otro aparece sólo una vez, en la presentación de los resultados ambos documentos aparecerán al mismo nivel, no pudiendo mostrar preferencia del uno sobre el otro.

Otra de las desventajas es que no se registra el contexto semántico de las palabras, incluso se pierde el orden en que aparecen las palabras en cada texto.

Este modelo se presume que en el cual se basa la implementación de Saber UCV y por eso es que en general se termina presentando el problema de que al usar el operador *AND* en las búsquedas exactas mencionadas en 1.5.2.1 se obtiene una alta **precision**¹ en los resultados pero un bajo **recall**² mientras que al usar el operador *OR* da una baja **precisión** y un gran **recall**.

¹Precision: la fracción, o porcentaje, de los documentos recuperados que son relevantes en la búsqueda efectuada.

²Recall: la fracción, o porcentaje, de los documentos relevantes en la colección que fueron recuperados por el sistema.

Con la propuesta del SSBSU se quiere una versión de recuperación de información que aplica métodos de mayor eficiencia, que permiten mejorar el desempeño (tiempo) y la calidad de los resultados.

2.5.2 Índices Invertidos:

Se denominan índices invertidos porque en vez de guardar los documentos con las palabras que en ellos aparecen, en estos se procede a guardar cada palabra y se indica los documentos en los cuales se encuentra y adicionalmente se puede registrar la posición en que aparece cada palabra con distintas granularidades, pudiendo ser estas: dentro del documento, del capítulo, del párrafo o de la oración. También pueden contener la frecuencia con que se presenta determinada palabra. Toda esta información nos permite mejorar los tiempos de búsqueda pero con ciertos costos.

El primero es el espacio en disco que implica guardar estos datos adicionales, que puede oscilar del 5% al 100% del valor inicial de almacenamiento, mientras que el segundo costo lo representa el esfuerzo computacional de actualizarlos una vez que se incorporan nuevos documentos (Mahapatra and Biswas, 2011).

Existen diversos tipos de **Índices Invertidos** y constantemente se están realizando investigaciones que permitan mejorar su desempeño motivado en que sobre ellos recae gran parte de la efectividad que podamos obtener ejecutando los *queries*. Algunos ejemplos de estos índices son el *Generalized Inverted Index* (GIN), también está el RUM³ o el VODKA⁴ que es otra implementación con menos literatura sobre posibles usos pero con métodos disponibles para su uso en manejadores de base de datos como PostgreSQL.

El espacio que ocupa la implementación de estos índices se puede ver afectado, por un lado tenemos que se puede reducir mediante el preprocesamiento que hagamos a las palabras buscando su raíz con el stemming {#steaming} o removiendo las stop words (las palabras que no generan mayor valor semántico como: la, el, tu).

Por otra parte el peso total se puede incrementar a medida que decidamos tener una granularidad más fina en el registro de las palabras y su ubicación dentro de los documentos. En el transcurso del desarrollo de nuestra investigación indicaremos en cuánto se incremento el espacio de almacenamiento en disco con la aplicación de este índice y la granularidad que se adoptó junto con el valor del costo en espacio de almacenamiento.

Continuando con los índices inversos hay estrategias que significan la adopción de generar dos índices inversos para un sistema, conteniendo uno de estos la lista de documentos y la frecuencia de la palabra, mientras que el otro registra la lista con las posiciones de la palabra.

El uso de los índices invertidos permite la denominada “búsqueda de texto completa” (*full text search*) que es uno de los pilares que sustenta a los motores de búsqueda y se entiende por este tipo de búsqueda aquella que permite encontrar documentos que contienen las palabras claves o frases determinadas en el texto del *query*. Adicionalmente podemos introducir el criterio de búsqueda de texto aproximado (*approximate text searching*), que permite flexibilizar la coincidencia entre el texto requerido y el resultado.

³En el vínculo <https://github.com/postgrespro/rum> se tiene acceso a la explicación e implementación de este índice para PostgreSQL.

⁴este índice fue presentado en la Postgres Conference en el año 2014 https://www.pgcon.org/2014/schedule/attachments/318_pgcon-2014-vodka.pdf

En la Solución que se propone, la optimización en la generación de este índice quedará bajo la administración del propio manejador de base de datos que es *postgreSQL*.

Cuando la base de datos que registra el índice invertido crece y no es viable almacenarla en un único computador, es necesario acudir al uso de técnicas que permitan distribuir la base de datos con el uso de tecnologías como *Spark*, *Hadoop*, *Apache Storm* entre otras.

2.5.3 Scoring Model:

Es un modelo aplicado en amplias colecciones de documentos donde es necesario exclusivamente mostrar al usuario que realizó la búsqueda, una fracción de los documentos encontrados, pero es necesario que estos sean los de mayor puntuación (*score*), de acuerdo a distintos criterios que permitan determinar cuales son los que tienen mayor relevancia, como lo puede ser la cantidad de veces que se repite una palabra aparecida en el *query* dentro del documento, o la distancia de aparición de cada una de las palabras que conforman el query (cantidad de palabras que median entre una y otra).

También se puede asignar un peso mayor en la generación del ranking, si una, o varias, de las palabras que generan el *query* aparece dentro del título, o en otros campos que se registrasen en la base de datos.

2.6 Procesamientos de texto:

En esta sección mostramos métodos de trabajo con los textos. Para el idioma español no son de abundante literatura, a diferencia de aquellos que están en el idioma inglés. Incluso hasta hace unos pocos años las herramientas computacionales para el procesamiento de los textos (*nlproc*) tampoco eran abundantes y sabiendo que son justamente los textos, el insumo que recibe de nuestro sistema de recuperación de información, la calidad en los procesamientos que sobre ellos hagamos, marcarán en gran medida la propia calidad del sistema que tengamos.

Como decíamos la literatura y herramientas disponibles para el NLP en el idioma español, fueron escasas durante un considerable período. Se tenían disponibles algunas como el *coreNLP* de la Universidad de Stanford pero no incluía todas las utilidades, tales como la identificación de parte del discurso (*Part of Speech Tagging*), ni el análisis morfológico (*Morphological Analysis*) o el reconocimiento de entidades nombradas (*Named Entity Recognition*), sino algunas pocas como el tokenizador 2.6.1.1 y el separador de oraciones (*Sentences Splitting*).

Casos similares se presentaban con otras herramientas, siendo un caso aparte el esfuerzo del CLiC-Centre de Llenguatge i Computació quienes hicieron la anotación del Corpus AnCora⁵. También la Universidad Politécnica de Cataluña creó la herramienta FreeLing⁶ que permitió realizar algunas de las funcionalidades mencionadas en el párrafo anterior. No obstante su integración en cadenas

⁵AnCora es un corpus del catalán (AnCora-CA) y del español (AnCora-ES) con diferentes niveles de anotación como lema y categoría morfológica, constituyentes y funciones sintácticas, estructura argumental y papeles temáticos, clase semántica verbal, tipo denotativo de los nombres deverbales, sentidos de WordNet nominales, entidades nombradas (NER), relaciones de correferencia (<http://cllc.ub.edu/corpus/es/ancora>)

< <http://cllc.ub.edu/corpus/es/ancora> >

⁶<https://nlp.lsi.upc.edu/freeling/node/1>

de trabajo y la actualización de sus modelos de entrenamiento, presentan rezagos en comparación a otros modelos que actualmente se están usando, los cuales serán evaluados en el transcurso de esta investigación y serán señalados en un capítulo que aún está por desarrollar.

2.6.1 Procesamiento del Lenguaje Natural (natural language processing- NLP):

Son las técnicas computacionales desarrolladas para permitir al computador “comprender” el significado de los textos de lenguaje natural. En esta Investigación son aplicadas una serie de estos métodos sobre los textos que componen el Corpus y podemos mencionar las siguientes:

2.6.1.1 Tokenizador:

Básicamente es separar el documento en palabras, o unidades semánticas que tengan algún significado a las cuales se le llaman *tokens*. Para el idioma español no representa un mayor reto, ya que se puede usar el espacio como delimitador de palabras, no así en otros idiomas como el chino donde el problema se aborda de manera distinta.

Al obtener las palabras como entidades separadas de un texto nos permite calcular la frecuencia de uso de las mismas.

Es común que las librerías de procesamiento de lenguaje natural contengan tokenizadores que presentan un 100% como métrica de precisión en el idioma español.

2.6.1.2 Entidades Nombradas (*named entity recognition-NER*):

Son los procesos que permiten la extracción de las distintas entidades contenidas dentro de los textos. Las entidades son: nombres, lugares, organizaciones. También permite detectar relaciones entre entidades. Las medidas de precisión en los módulos de NER alcanzan una medida cercana al 89% en modelos entrenados con *machine learning*, tal es el caso de spacy, que es una de las librerías propuestas para realizar estos procesamiento en nuestro desarrollo.

2.6.1.3 Etiquetado de Partes del Discurso (*Part of speech tagging-POS*):

Una de las técnicas usadas en el Procesamiento del Lenguaje Natural es el part-of-speech (POS) y consiste en asignar un rol sintáctico a cada palabra dentro de una frase (Eisenstein, 2019) siendo necesario para ello evaluar cómo cada palabra se relaciona con las otras que están contenidas en una oración y así se revela la estructura sintáctica.

Los roles sintácticos principales de interés en la elaboración de esta Investigación son los sustantivos, adjetivos y verbos.

- Los sustantivos tienden a describir entidades y conceptos.
- Los verbos generalmente señalan eventos y acciones.
- Los adjetivos describen propiedades de las entidades

Igualmente dentro del POS se identifican otros roles sintácticos como los adverbios, nombres propios, interjecciones entre otros.

El POS es un procesamiento que sirve de insumo para la coocurrencia de palabras, que es una de las formas en que se representan los resultados de los *querys* en la SSCSU.

En el estado del arte este etiquetado alcanza un 98% de precisión.

2.6.1.4 Stemming:

Es el proceso en que se consigue el lema de una palabra, entendiendo que el lema es la forma que por convenio se acepta como representante de todas las formas flexionadas de una misma palabra.

Al buscar el lema se tiene presente la función sintáctica que tiene la palabra, es decir que se evalúa el contexto en el que ocurre. Una de las ventajas de aplicar esta técnica es que se reduce el vocabulario del Corpus y eso conlleva a que también se reduce el espacio de búsqueda en los documentos.

En el estado del arte este etiquetado alcanza un 96% de precisión.

2.7 Minería de Texto:

La extracción de ideas útiles derivadas de textos mediante la aplicación de algoritmos estadísticos y computacionales, se conoce con el nombre de minería de texto, analítica de texto o aprendizaje automático para textos (text mining, text analytics, machine learning from text). Se quiere con ella representar el conocimiento en una forma más abstracta y así poder detectar relaciones en los textos.

El uso de las técnicas de minería de texto ha ganado atención en recientes años motivado a las grandes cantidades de textos digitales que están disponibles. La minería de texto surge para dar respuesta a la necesidad de tener métodos y algoritmos que permitan procesar estos datos no estructurados (Aggarwal and Zhai, 2012b). Una vez que mediante el proceso de recuperación de información tenemos un cúmulo de textos, es posible que necesitemos dar un paso más allá y usar un conjunto de técnicas que nos permitan analizar y digerir estos textos, mediante la detección de patrones.

Uno de los desafíos al trabajar con textos es darles estructura para que resulte viable trabajar con ellos desde la perspectiva de procesos computacionales.

Una de las primeras fases consiste en agrupar todos los textos en un Corpus. Posteriormente se procederá a conformar una matriz dispersa de una alta dimensionalidad que se denominará “*Sparse Term-Document Matrix*”) de tamaño $n \times d$, donde n es el número total de documentos y d es la cantidad de términos o vocabulario (palabras distintas) presentes entre todos los documentos. Formalmente sabemos que la entrada (i,j) de nuestra matriz es la frecuencia (cantidad de veces que aparece) de la palabra j en el documento i .

El problema de la alta dimensionalidad de la matriz mencionada motiva ir aplicando otras técnicas que en principio puedan colaborar a reducirla por medio de simplificar los atributos, es decir, disminuyendo el vocabulario por ejemplo aplicando el algoritmo de Porter (stemming).

Igualmente a nivel de minería de texto se hace deseable poder contar con la identificación semántica de cada una de las palabras que conforman nuestro vocabulario, para así obtener

2.7. MINERÍA DE TEXTO:

representaciones que aportan un mayor significado. Los procesamientos inherentes al NLP mencionados anteriormente son insumo para la minería de texto.

2.7.1 Coocurrencia de Palabras:

En esta investigación se usará una técnica denominada “Coocurrencia de Palabras” para la detección de patrones en los textos. Esto consiste en evaluar las palabras que coocurren dentro de los documentos que conforman el *corpus* y se puede hacer con distintas granularidades. Por ejemplo: las palabras que coocurren una seguida de otra o las que coocurren dentro de la misma oración, o dentro de un párrafo y así sucesivamente.

Para la representación visual se usan los grafos representando cada palabra un nodo y la coocurrencia de una palabra con otra implica que se extienda un arco entre ellas. Las palabras dispuestas para representarse en el grafo serán exclusivamente las que tengan la función dentro del discurso (POS) 2.6.1.3 de adjetivos y sustantivos, es decir que cada coocurrencia será un sustantivo con el adjetivo que la acompaña, donde es posible tener una relación de un sustantivo con $\{0,1,\dots,n\}$ adjetivos.

La selección de las funciones gramaticales propuestas se hace para disminuir el espacio de representación y se considera que los sustantivos, al contar con el adjetivo que las acompaña, logran hacer una representación que muestra proximidad semántica y se representan los temas (*topics*) más relevantes (Segev, 2021).

En el método que se usará en este Sistema se filtrarán las n (n igual a 100), palabras que presenten mayor coocurrencia dentro de los resúmenes filtrados en el *query*, siendo posible seleccionar la granularidad (todo el texto o en un párrafo).

En la figura 2.2 podemos visualizar lo expuesto de una manera gráfica al ver la representación en un grafo de la coocurrencia de palabras sobre los textos de los resúmenes de las Tesis y TEG de la Escuela de Física de la U.C.V.

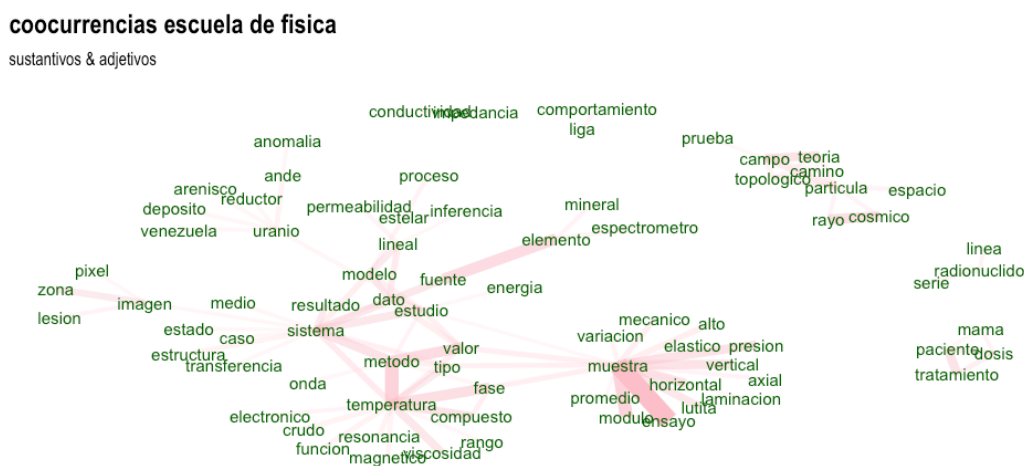


Figura 2.2: Coocurrencia de Palabras

2.8 Similitud de documentos:

Para poder realizar la recomendación de documentos, una de las técnicas que se usa es medir la similitud que presenta un documento con los otros contenidos en el corpus. Un ejemplo de esta técnica es el uso de la similitud coseno que se explica con esta fórmula.

$$\cos(\mathbf{t}, \mathbf{e}) = \frac{\mathbf{t} \cdot \mathbf{e}}{\|\mathbf{t}\| \|\mathbf{e}\|} = \frac{\sum_{i=1}^n \mathbf{t}_i \mathbf{e}_i}{\sqrt{\sum_{i=1}^n (\mathbf{t}_i)^2} \sqrt{\sum_{i=1}^n (\mathbf{e}_i)^2}} \quad (2.1)$$

En la fórmula t representa un documento y e representa otro documento. Ambos documentos se asumen que están en un espacio con i atributos, o dimensiones, y la intención es calcular un índice de similitud entre ambos documentos.

Este es uno de los métodos más usados para detectar similitudes en los textos, aunque existen otras fórmulas para el cálculo de la similitud como es el índice de jaccard.

El otro elemento de gran importancia en obtención de esta medición, es la representación que se haga del documento. Son distintas las técnicas que existen estando entre ellas la representación mediante “bolsas de palabras” o *bag of words*. Recientemente se han creado formas más complejas para la representación como lo son los *words embeddings* que son obtenidos mediante el entrenamiento de redes neuronales de aprendizaje profundo.

Estas formas de representación de textos serán investigadas más a fondo en futuras fases de la elaboración de este trabajo.

2.9 Sistemas Distribuidos:

Los distintos procesos y componentes de la Solución propuesta han sido diseñados e implementados como un sistema distribuido y por eso haremos una mención a este tema.

Una definición formal que se le puede dar a los sistemas distribuidos es “cuando los componentes de hardware y/o software se encuentran localizados en una red de computadores y estos coordinan sus acciones sólo mediante el pase de mensajes” (Coulouris, 2012).

Algunas de las principales características que tienen los sistemas distribuidos es la tolerancia a fallos, compartir recursos, concurrencia, ser escalables (Czaja, 2018) entre otras. Mencionamos estas en particular al ser propiedades que están presentes en la propuesta acá descrita.

1. Fiabilidad o la llamada tolerancia a fallos: en caso de fallar un componente del sistema los otros se deben mantener en funcionamiento.
2. Compartir recursos: un conjunto de usuarios pueden compartir recursos como archivos o base de datos.
3. Concurrencia: poder ejecutar varios trabajos en simultáneo.
4. Escalable: al ser incrementada la escala del sistema se debe mantener en funcionamiento el sistema sin mayores contratiempos.

2.9.1 Contenedores:

Un contenedor es una abstracción de una aplicación que se crea en un ambiente virtual, en el cual se encuentran “empaquetados” todos los componentes (sistema operativo, librerías, dependencias, etc.), que una aplicación necesita para poder ejecutarse. En su diseño se tiene presente que sean ligeros y que con otros contenedores pueden compartir el *kernel*, usando un sistema de múltiples capas, que también pueden ser compartidas entre diversos contenedores, ahorrando espacio en disco del *host* donde se alojan los contenedores.

El uso de los contenedores permite crear, distribuir y colocar en producción aplicaciones de software de una forma sencilla, segura y reproducible. También a cada contenedor se le puede realizar una asignación de recursos (memoria, cpu, almacenamiento) que garantice un óptimo funcionamiento de la aplicación que contienen.

Es importante señalar que el uso de esta tecnología añade un entorno de seguridad al estar cada contenedor en un ambiente aislado.

Para cada contenedor es necesario usar una imagen donde se definen las dependencias necesarias para su funcionamiento.

2.9.2 Orquestador:

Al tener diversos contenedores, donde cada uno aloja una aplicación distinta, puede resultar necesario que todos se integren en un sistema. Para que esta integración sea viable es necesario contar con un orquestador. Su uso permitirá lograr altos grados de portabilidad y reproducibilidad, pudiendo colocarlos en la nube o en centros de datos garantizando que se pueda hacer el *deploy* de forma sencilla y fiel a lo que se implementó en el ambiente de desarrollo.

En el caso de la Solución propuesta se adoptará el uso de Docker Compose como orquestador y en el Capítulo que contiene la Propuesta Técnica 3.1 serán expuestas las funcionalidades de cada contenedor y se apreciará la integración que brindará el orquestador.

2.10 Tendencias actuales Sistemas de Información:

Si bien anteriormente las búsquedas de información dentro de un conjunto de textos se procesaban determinando la aparición o no de palabras, o de frases dentro de un determinado texto, este método ha ido evolucionando para llegar hoy en día a un elevado nivel de abstracción, donde a partir de la necesidad de obtener una determinada información, es decir, de aquello que necesitamos buscar, que antes consistía en hacer *match* con un objeto de información, hemos pasado de los motores de búsqueda (*search engines*) a los motores de respuestas (*answering engines*) (Balog, 2018), donde el sistema ante un determinada consulta del usuario va a retornar una serie de resultados enriquecidos, mostrando la identificación de entidades, hechos y cualquier otro dato estructurado que esté de forma explícita, o incluso implícita, mencionado dentro de los textos que conforman el corpus.

Para lograr la generación de estos resultados se han tenido que conformar las llamadas bases de conocimiento o *knowledge bases*, que son repositorios donde previo a la búsqueda de la

información dentro del sistema, se logra ir estructurando la organización de la información alrededor de objetos o datos específicos que se denominan **entidades**. Estos conceptos y métodos se asocian directamente a los que también se proponen, de manera más amplia, en la denominada *web semántica*⁷.

Como ejemplo de una *knowledge bases* se puede mencionar a DBpedia, que se encuentra en la dirección www.dbpedia.org y es un proyecto en donde puede acceder a una red global y unificada de grafos de conocimiento, la cual cubre más de 20 idiomas y principalmente genera sus grafos de conocimiento a partir de las publicaciones del Proyecto Wikipedia.

2.10.0.1 Temas en proceso de investigación:

Por estar aún en curso esta investigación a continuación se mencionan algunos temas que se está evaluando incluir en este marco teórico

- Métodos usados para el almacenamiento o la indexación como crear agrupamientos (*clusters*) de aquellos documentos que compartan algunas características, por ejemplo, en la temática que aborden. Otra de las innovaciones que se están añadiendo a los sistemas de recuperación de información, es generar resúmenes con técnicas de procesamiento de lenguaje natural soportadas en el uso de arquitecturas de aprendizaje profundo (*deep learning*).
- Resultados ante una búsqueda personalizados: al existir mecanismos como la sesión de usuario o las *cookies* que guardan información contextual, permitiendo que ante un mismo *query* en distintos equipos, los resultados sean distintos en función de la persona que hizo la búsqueda.
- En cuanto al estado del arte existen distintas librerías y modelos de representación de documentos, palabras y caracteres. Algunos de estos modelos son fastext, word2vec, GPT3. Para este momento aún estamos haciendo la evaluación del uso de ellos para nuestra investigación, por lo cual, sólo los mencionamos referencialmente en este Anteproyecto.

⁷Se basa en la idea de añadir metadatos semánticos y ontológicos a la *World Wide Web*. Esas informaciones adicionales —que describen el contenido, el significado y la relación de los datos— se deben proporcionar de manera formal, para que así sea posible evaluarlas automáticamente por máquinas de procesamiento. Tomado de Wikipedia

Capítulo 3

Propuesta Técnica:

En este Capítulo se detalla la propuesta técnica de la Solución en la sección 3.1. Posteriormente se indica el Objetivo General 3.2 y los Objetivos Específicos 3.3 que se aspiran cubrir en el desarrollo. A continuación se muestra el esquema general de la aplicación 3.4, junto con el esquema de funcionamiento en los contenedores 3.4.2 y se describe el funcionamiento de cada contenedor.

Se realiza un análisis de factibilidad en la sección 3.5 para la implementación y se muestra un cronograma 3.6 de actividades.

3.1 Propuesta Técnica:

Crear una Solución que funcione como una aplicación web distribuida bajo la arquitectura cliente-servidor. Del lado del servidor se encontrarán contenedores que alojarán los distintos servicios necesarios para el funcionamiento.

Del lado del cliente se podrán formular los *queries* con múltiples atributos (la frase a buscar, la jerarquía académica, intervalo de fechas, facultad y/o escuela de adscripción).

La Solución permitirá generar procesos de extracción de información (*information retrieval*) y los resultados se mostrarán en tablas interactivas, gráficos y grafos de coocurrencia de palabras. Adicionalmente se indicarán recomendaciones de textos basados en la similitud que presente un documento con los otros.

Los textos del resumen de cada tesis con los cuales se conformará el *Corpus*, serán obtenidos y actualizados mediante técnicas de *web crawling* realizadas al repositorio Saber UCV.

Se opta por la estrategia de hacer el *web crawling* a los documentos contenidos en Saber UCV al no ser posible, al momento de realizar esta propuesta, la obtención de la base de datos de los documentos ahí alojados.

La Solución propuesta está diseñada para irse actualizando periódicamente incorporando los nuevos documentos que sean alojados en el repositorio Saber UCV.

Contará con un procedimiento que permite clasificar los documentos (tesis) por área académica, según donde haya efectuado estudios el autor del trabajo, al ser la tesis el requisito necesario para obtener el título de grado. Actualmente el repositorio Saber UCV no dispone de esta clasificación.

3.2 Objetivo:

Crear una Solución que permita realizar la clasificación de documentos y la búsqueda de información sobre los trabajos de investigación que se encuentran en el Repositorio SABER UCV usando técnicas de extracción de información (*information retrieval*).

Para el procesamiento de las búsquedas será usado un sistema distribuido con distintos componentes (contenedores) que permitan la ingesta, procesamiento y transformación de los datos, al igual que el alojamiento de los mismos en una base de datos.

3.3 Objetivos Específicos

- *Querys* multi atributo en distintas dimensiones: tiempo, jerarquías (pregrado, especializaciones, maestría y/o doctorado), facultad ,carreras o postgrados.
- A cada resultado que arroje la búsqueda se le debe asignar un nivel de relevancia. El conjunto de los textos recuperados se debe mostrar ordenadamente de mayor a menor relevancia. La asignación de la relevancia será hecha mediante una función que cuente con distintos parámetros.
- Los datos: textos, fechas, clasificaciones, y demás que sean necesarios para el funcionamiento de la Solución, deben registrarse en una base de datos estructurada y la ingesta y consulta se efectuará mediante un manejador de base de datos.
- Generar recomendaciones de textos basados en la similitud coseno que presente un documento con cada uno de otros que conforman el *Corpus*.
- En los documentos extraídos se debe contar con el enlace a los documentos que reposan en Saber UCV.
- Permitir la concurrencia de accesos al Sistema.
- La tolerancia a fallas en los contenedores.
- Contar con el certificado SSL para acceso seguro por parte de los visitantes.
- Procesar las coocurrencias de las palabras más comunes las cuales se deberán visualizar mediante grafos.
- Disponer de una aplicación web para el acceso al Sistema por parte del cliente.
- Disponer de ventanas emergentes de ayuda en la interface gráfica de la aplicación web.
- En la representación de coocurrencias de palabras mediante grafos, se debe poder filtrar documentos interactivamente al hacer *click* con el *mouse* sobre los arcos que unen un par de nodos, donde cada uno de estos representa la coocurrencia de una dupla de palabras. El *click* generará el evento para el query sobre los documentos que contienen esa determinada coocurrencia.

3.4 Esquema del SSCSU:

Podemos representar el Sistema y sus interacciones mediante un esquema con un elevado grado de abstracción. Ver figura 3.1.

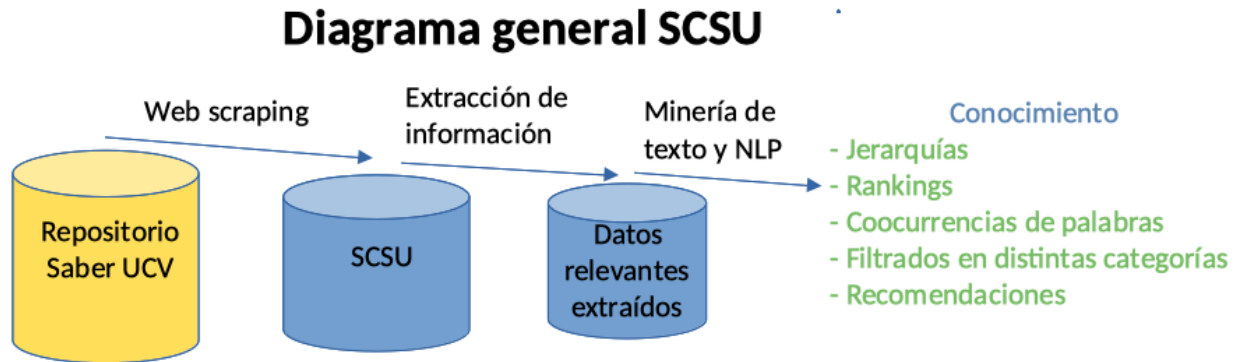


Figura 3.1: Esquema General

En la figura 3.1 se muestra como el SSCSU está diseñado para que mediante técnicas de *web crawling* se pueda extraer la información del repositorio Saber UCV. Posteriormente el sistema ante la consulta del usuario va a recuperar los documentos que sean relevantes, representando el conocimiento con la generación de jerarquías, rankings, coocurrencias de las palabras más mencionadas y recomendaciones de texto según similitudes.

El SSCSU se implementará mediante el uso de Docker. En la figura 3.2 se muestra el esquema con los contenedores.

3.4.1 Cliente - Servidor:

La arquitectura **cliente-servidor** en esta Propuesta se basa en:

3.4.1.1 1 - Cliente - Navegador web:

El acceso del cliente a la aplicación web se hace en el navegador web. Este realiza la petición al servidor que aloja la Solución. Actualmente se cuenta con un prototipo al que se puede acceder en la dirección <https://proyecto.me/>.

El sistema no está diseñado para usarse desde dispositivos móviles, aunque igualmente es viable el acceso no estando optimizada la interfaz del usuario ni las visualizaciones de los resultados de los *queries*.

3.4.1.2 2 - Servidor:

El prototipo del sistema requiere para funcionar al menos estos recursos:

- 2 CPU virtual

DIAGRAMA COMPONENTES DOCKER

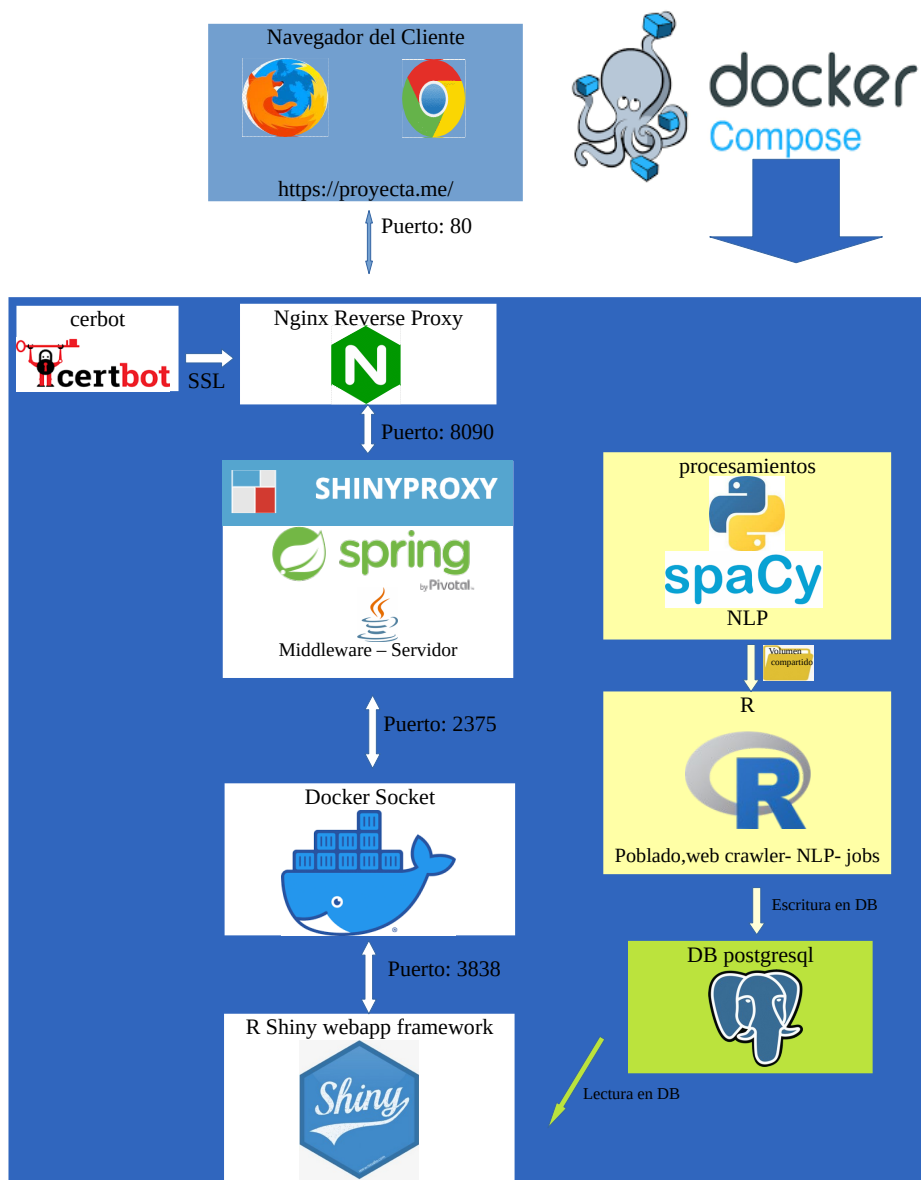


Figura 3.2: Diagrama de Componentes

3.4. ESQUEMA DEL SSCSU:

- 2 GB de memoria RAM
- 50 GB de disco duro

En el servidor se instala el software Docker sobre el cual se despliegan los siguientes contenedores:

3.4.2 Contenedores:

3.4.2.1 Docker Compose:

Funciona como un orquestador para correr aplicaciones distribuidas en múltiples contenedores usando un archivo en formato yml donde se establecen las imágenes, los puertos y los volúmenes que serán usados y compartidos por cada uno de los contenedores.

3.4.2.2 Contenedor con Nginx:

Es el servidor web/proxy inverso de código abierto que sirve para redireccionar las peticiones del puerto 80 al puerto 8090. Mediante este contenedor también se define el certificado SSL para permitir conexiones por el protocolo HTTPS.

Este contenedor fue generado desde una imagen oficial de NGINX que se encuentra en el ***docker hub***¹ sin añadir ninguna capa (layer) adicional.

3.4.2.3 Contenedor con Cerbot:

Cerbot es una herramienta de código abierto que permite habilitar las conexiones mediante el protocolo HTTPS con el uso de un certificado “Let’s Encrypt”. El uso de este certificado está asociado al uso de un dominio en el *deploy* de la aplicación.

Este contenedor fue generado desde una imagen de CERBOT del ***docker hub*** sin ninguna modificación posterior.

3.4.2.4 Contenedor con Shinyproxy:

Es una implementación del servidor “*Spring boot*” que dará servicio a las aplicación web desarrolladas en *shiny*² 3.4.2.5. Con el uso de este *middleware* se obtienen las siguientes ventajas:

- Ante cada petición de acceso al servidor se despliega un *workspace* completamente aislado, es decir, un contenedor distinto. Las aplicaciones desarrolladas en *shiny* son *single threaded* y adoptar esta estrategia representa una ventaja, motivado a que se pueden controlar los recursos de memoria y cpu asignados a cada contenedor que se despliega.
- Permite establecer *login* en el uso de la aplicación y grupos de usuarios. También da soporte a distintos métodos de autenticación. Si bien en estos momentos la aplicación es de libre acceso, en algún momento se pudiese restringir y no sería necesaria ninguna modificación en la arquitectura, más allá de cambiar el archivo de configuración.

¹repositorio de imágenes de contenedores de docker.

²Shiny un framework para crear aplicaciones web en el lenguaje de programación R.

- Uso de una tecnología estable y probada.

Es necesario destacar que originalmente *Shiny* como *framework* cuenta con su propio software que actúa como servidor, pero para tener acceso a ciertas funcionalidades es necesario pagar por el servicio directamente a la Fundación RStudio Software y usar el alojamiento que ellos proveen, no siendo todos los componentes de código abierto. Por ejemplo, el acceso mediante *login* no está disponible en la versión libre del *Shiny Server* sino en la *Shiny Server Professional*.

Ciertas configuraciones de librerías e incluso la propia contenerización de la aplicación, no es posible usando el servicio pago, así que la propuesta acá adoptada, si bien representa un mayor esfuerzo en la configuración, claramente implica que se obtienen una serie de ventajas, por todas las adaptaciones y control que es posible realizar al y sobre el sistema.

Este contenedor funciona en el puerto 2375 y fue generado desde la imagen del docker hub *Shinyproxy* sin ninguna modificación.

3.4.2.5 Contenedor con “R Shiny Web App framework”:

En este contenedor es donde reposa la aplicación web con todas las librerías necesarias para generar las visualizaciones y remitir los *queries* al contenedor del manejador de la base de datos 3.4.2.6 . Como comentamos anteriormente, cada vez que ocurre desde el navegador del cliente una petición de acceso, desde el contenedor *shinyproxy* 3.4.2.4, se crea una replica de este contenedor con todos los elementos necesarios para que la app funcione correctamente.

En caso de presentar alguna falla, el sistema sería tolerante, porque se pueden seguir recibiendo peticiones que replicarían una imagen nueva del contenedor sin afectar al que presentase el fallo, o viceversa.

Desde este contenedor se realiza el acceso de lectura al contenedor que contiene *PostgreSQL* 3.4.2.6 donde reposa la base de datos que contiene los textos ya procesados.

Por los momentos no hay escritura de datos en las tablas, pero está contemplado que se registren los *queries* formulados en alguna tabla, junto con los documentos que el usuario revisa mediante las interacciones, para así generar métricas de calidad del sistema y del uso que se le da.

La imagen que se usa en este servidor fue definida a medida con todos los recursos necesarios.

En un posterior Capítulo a desarrollar en donde se mostrará el proceso de desarrollo de la Solución, serán descritas todas las librerías que se encuentran incluidas en este contenedor y se expondrán las razones para seleccionar cada una de ellas.

Varios de los procesos que se ejecutan en este contenedor ocurren al momento de recibir un *query*, no obstante todos los procesos que puedan ser pre computados, se trata de ejecutarlos previamente en el contenedor “R Servicios” 3.4.2.7, para lograr así la disminución de los tiempos.

Funcionalidades principales: **ENTRADAS**

- Contiene un campo para la entrada de texto que generará el query.
- Contiene un selector para indicar si se quiere generar la coocurrencia de palabras
- Contiene tablas para seleccionar:

3.4. ESQUEMA DEL SSCSU:

- 1) Nivel académico del trabajo. Opciones (pregrado, especialización, maestría, doctorado).
- 2) Facultad o Centro de adscripción. Opciones: 11 Facultades más un centro (CENDES).
- 3) Nombre del pregrado o postgrado. En total son 412 las opciones.

Cada una de las tablas anteriores se actualiza según se vayan seleccionando las relaciones y la disponibilidades. Por ejemplo, al seleccionar pregrado solo se mostrarán los nombres de las carreras de pregrado, pero si se selecciona también el nombre de la Facultad, sólo se mostrarán las carreras de pregrado dentro de la Facultad seleccionada. Para una determinada tabla también se permiten selecciones múltiples dando una total flexibilidad al momento de ejecutar los *queries*.

SALIDAS - Ante el *query* se genera:

- Una tabla creada con la librería *reactable*. En la misma se muestra cada uno de los documentos recuperados con los distintos atributos disponibles: autor, fecha, palabras claves, texto resumen. El orden en que aparece está generado con una función de ranking. Adicionalmente se muestra un enlace al repositorio Saber UCV que es donde se encuentra alojado el respectivo trabajo (el documento en PDF). Igualmente se presentan los textos que tienen mayor similitud con el texto seleccionado.
- Un gráfico con la frecuencia por año de los trabajos extraídos mediante el *query*. El gráfico se generará con la librería *apexcharter* que es un wrapper para Javascript, por lo cual el gráfico tiene ciertas interactividades. Con el *hover* muestra el valor de cada columna.
- Gráfico de coocurrencia interactivo de palabras: se genera mediante la librería de *VisNetwork* que también es un *wrapper* de javascript. Este gráfico permite seleccionar un arco de unión entre dos palabras coocurrentes. Al realizar la selección se filtran un subconjunto de los documentos que contienen ambas palabras representadas por nodos. Los documentos filtrados se mostrarán en una tabla contigua, también generada en *reactable*, donde sólo se incluye el texto resumen de cada trabajo. En un próximo Capítulo a desarrollar será mencionado el diseño y funcionamiento a detalle de esta visualización.
- Gráfico de coocurrencia estático de palabras: mediante la librería *ggraph* son generados un par de gráficos con distintas granularidades. El primero exhibe la misma coocurrencia de palabras expuesta en el punto anterior pero esta vez es generada en una visualización estática. En cuando a la granularidad se muestran las palabras que coocurren dentro de todo el resumen, El segundo gráfico también muestra la coocurrencia, pero solo de palabras que coocurren una seguida de otra dentro del texto resumen, es decir que se muestran los resultados con una menor granularidad.

Con la librería *UdPipe* se generan las estructuras de datos necesarias para generar los grafos (arcos y nodos).

3.4.2.6 Contenedor con PostgreSQL:

En este contenedor tenemos una imagen de *PostgreSQL* versión 13.3. No fue realizada ninguna otra modificación distinta a la definición de usuarios y poblado desde base de datos incluyendo

un volumen compartido para garantizar que tengamos “datos persistentes. Este contenedor recibe consultas del contenedor”R Shiny Web App framework” 3.4.2.5 y escritura desde el contenedor “R imagen Servicios” 3.4.2.7.

En este contenedor ocurre la indexación de la base de datos y la generación del ranking al procesar el resultado con la función de PostgreSQL llamada *tsrank* ³.

En una tabla se encuentra la identificación del documento, la fecha de creación y propiamente los textos que mediante el *Tsvector* ⁴ almacena el título, el resumen, el autor y las palabras claves.

En otra tabla se encuentra el almacenamiento del procesamiento que se le hace a los textos, clasificando cada una de las palabras mediante el *part of speech*, y registrando el identificador del documento al que está asociada y el lema ⁵ de cada una.

El *TSvector* es la estructura de datos que permite la búsqueda de texto completa (*Full Text Search*) mediante la función de *PostgreSQL* denominada *tsquery*. En un futuro capítulo a desarrollar será mostrado el diseño de las tablas con sus campos.

3.4.2.7 Contenedor con “R Imagen Servicios”:

En este contenedor se creó una imagen con todos los servicios necesarios para realizar el web crawling, el procesamiento de textos y la descarga de los archivos desde Saber UCV para realizar la clasificación de las Tesis. Al iniciar el Sistema también contiene las funcionalidades que permiten realizar la creación de la base de datos, de las tablas y el poblado de estas.

Periodicamente es invocado este contenedor para realizar los procesos de incorporación de aquellos documentos nuevos que se detecte que están disponibles en Saber UCV.

En una primera fase se usó la utilidad del sistema operativo linux para la programación de actividades, sin embargo se está implementando el uso de la tecnología *Apache Airflow* para la ejecución de la programación de los flujos de trabajo.

La imagen base usada es la del proyecto ***Rocker*** (Boettiger and Eddelbuettel, 2017), la cual es una versión ampliamente probada y optimizada en, y por, la comunidad de usuarios de R.

Posteriormente serán descritas todas las librerías que fueron añadidas mediante una capa (layer) a este contenedor. Por los momentos podemos citar que se encuentran agregadas las siguientes:

3.4.2.7.1 Text Mining y NLP: Generalmente las distintas librerías que permiten realizar procesos de *Natural Language Processing* también hacen procesos de *Text Mining* parcial o totalmente. En la investigación fueron evaluadas múltiples librerías como *Spacy*, *Quanteda*, *OPENLP*, *CoreNLP*, *Freeling* y *Udpipe*.

En una futura entrega se hará una breve mención a cada una y la razón por la cual se adoptó *Spacy* para varios de los procesos de NLP.

³ esta función mide la relevancia de los documentos para una consulta en particular, de modo que cuando haya muchas coincidencias, las más relevantes puedan mostrarse primero tomando en cuenta la información léxica, de proximidad y estructural del documento (título, cuerpo del documento, etc).

⁴ El *Tsvector* es una función de postgresql que crea una estructura de datos que es una lista ordenada de distintos lexemas, que son palabras que se han normalizado para fusionar diferentes variantes de la misma palabra mediante el algoritmo de Porter.

⁵ El lema es la forma que por convenio se acepta como representante de todas las formas flexionadas de una misma palabra.

3.5. FACTIBILIDAD:

Para ejecutar *Spacy* es necesario usar los archivos que se encuentran en el contenedor “*Python Spacy*” , 3.4.2.8 donde está instalada la librería *Spacy* que corre en *Python*.

Procesos que se ejecutan llamando al contenedor “*Python Spacy*”: Tokenización, POS y Lematización.

También en este contenedor se realizan los siguientes procesos:

1. Poblado base de datos: se hace el *web crawling* para el poblado inicial y adicionales de la base de dato, con los textos de los Resúmenes.
2. Clasificación de los documentos: mediante una rutina compuesta por varios algoritmos serán clasificados los distintos trabajos según lo antes expuesto.
3. Cálculo de la Similitud de los documentos: cuando se ha realizado la lematización de las palabras se procede a generar una matriz de tipo Td-Idf (term document- inverse document frequency), que sirve de insumo para el cálculo de similitud entre los documentos. Este cálculo de similitud se realiza con la librería *Quanteda.textstats* y se usa la medida de similitud coseno, ya que varios autores la sitúan como una de las mejores formas de comparar la similitud entre un documento y otro.

3.4.2.8 Contenedor con “*Python Spacy*”:

Se creó una imagen que contiene un ubuntu con python, spacy y el modelo de Spacy es_core_news_sm . Su función es que mediante un volumen compartido pueda ser invocado desde el contenedor “*R Imagen Servicios*” 3.4.2.7 para así realizar los procesamientos de NLP antes descritos.

3.5 Factibilidad:

Para la propuesta del SCSU se hizo una evaluación de la factibilidad del desarrollo del proyecto que consistió en hacer pruebas de *web crawling* sobre el repositorio Saber UCV al no ser viable la obtención del conjunto de datos por otro medio. Igualmente se realizaron pruebas sobre el hardware disponible con arquitecturas similares a la que se propondrá más adelante en nuestro desarrollo.

Estas pruebas fueron exitosas por lo cual no se estima que exista algún factor que impida la implementación del Sistema acá expuesto.

3.6 Cronograma:

El cronograma de desarrollo propuesto es el siguiente:

Cuadro 3.1: Cronograma de desarrollo de la Solución

ID Tarea	cantidad medida	predecesora
1 web craling inicial	2 semanas	0
2 diseño algoritmo clasificación	3 semanas	1
3 diseño componentes de la aplicación	4 semanas	1
4 Implementación de componentes	4 semanas	3
5 Pruebas contenedores Docker	1 semanas	4
6 Pruebas de software	1 semanas	5
7 implementación en servidor	1 semanas	6
8 Elaboración de trabajo especial de grado	8 semanas	7
TOTAL	24 semanas	

Bibliografía

- Aggarwal, C. C. (2018). *Machine Learning for Text*. Springer International Publishing : Imprint: Springer, Cham, 1st ed. 2018 edition. DOI: 10.1007/978-3-319-73531-3.
- Aggarwal, C. C. and Zhai, C., editors (2012a). *Mining text data*. Springer, New York Heidelberg. DOI: 10.1007/978-1-4614-3223-4.
- Aggarwal, C. C. and Zhai, C., editors (2012b). *Mining text data*. Springer, New York Heidelberg. DOI: 10.1007/978-1-4614-3223-4.
- Balog, K. (2018). *Entity-Oriented Search*. Number 39 in The Information Retrieval Series. Springer International Publishing : Imprint: Springer, Cham, 1st ed. 2018 edition. DOI: 10.1007/978-3-319-93935-3.
- Boettiger, C. and Eddelbuettel, D. (2017). An introduction to rocker: Docker containers for r. *The R Journal*, 9(2):527–536.
- Coulouris, G. F., editor (2012). *Distributed systems: concepts and design*. Addison-Wesley, Boston, 5th ed edition.
- Czaja, L. (2018). *Introduction to Distributed Computer Systems: Principles and Features*. Number 27 in Lecture Notes in Networks and Systems. Springer International Publishing : Imprint: Springer, Cham, 1st ed. 2018 edition. DOI: 10.1007/978-3-319-72023-4.
- Dueñas, M., Rojas, D., and Morales, M. (2011). Propuesta metodológica para realizar mapas de conocimiento. *Revista Facultad de Ciencias Económicas*, 20(1):77–90.
- Eisenstein, J. (2019). *Introduction to natural language processing*. Adaptive computation and machine learning. The MIT Press, Cambridge, Massachusetts.
- Fredkin, E. (1960). Trie memory. *Communications of the ACM*, 3(9):490–499.
- Goodrich, M. T., Tamassia, R., and Goldwasser, M. H. (2013). *Data structures and algorithms in Python*. Wiley, Hoboken, NJ.
- Guevara, J. (2015). *Desarrollo de un Motor de Búsquedas para la recuperación de documentos académicos de la Facultad de Ciencias*. PhD thesis. Accepted: 2015-05-25T14:29:29Z.
- Hernández Orallo, J., Ramírez Quintana, M., and Ferri Ramírez, C. (2004). *Introducción a la minería de datos*. Pearson Educación, Madrid. OCLC: 989816167.
- Inmon, W. H. (2002). *Building the data warehouse*. J. Wiley, New York, 3rd ed edition.

- Knuth, D. E. (1997). *The art of computer programming*. Addison-Wesley, Reading, Mass, 3rd ed edition.
- Kraft, D. H. and Colvin, E. (2017). Fuzzy information retrieval. *Synthesis Lectures on Information Concepts, Retrieval, and Services*, 9(1):i–63.
- Lehman, R. (2007). Learning object repositories. *New Directions for Adult and Continuing Education*, 2007(113):57–66.
- Mahapatra, A. K. and Biswas, S. (2011). Inverted indexes: Types and techniques. *International Journal of Computer Science Issues*, 8.
- Manning, C. D., Raghavan, P., and Schütze, H. (2008). *Introduction to information retrieval*. Cambridge University Press, New York. OCLC: ocn190786122.
- Rodríguez Laguna, J. (2016). *Sistema de recomendación para el buscador académico venezolano*. PhD thesis. Accepted: 2016-03-16T17:26:20Z.
- Segev, E. (2021). *Semantic Network Analysis in Social Sciences*. Routledge.
- Sánchez, J. (2008). *Elaboración de un prototipo de buscador de documentos académicos de la Facultad de Ciencias*. PhD thesis, Caracas, Venezuela.
- Zhai, C. and Massung, S. (2016). *Text data management and analysis: a practical introduction to information retrieval and text mining*. Number #12 in ACM Books. ACM Books, New York, first edition edition. OCLC: ocn957355971.
- Zhang, J. (2008). *Visualization for information retrieval*. The information retrieval series. Springer, Berlin. OCLC: ocn173239487.