

Motor de Búsqueda usando Hadoop y algoritmos: Índice Invertido y Pagerank

Patrick Xavier Marquez Choque*
Universidad Católica San Pablo

Jean Carlo Cornejo Cornejo†
Universidad Católica San Pablo

Mayo 2022

1 Desarrollo

Para la implementación de dicho trabajo, se nos dió la tarea de implementar un clúster en *Hadoop* y implementar el algoritmo de índice invertido y *page rank* para lograr un motor de búsqueda de términos de data dentro de un sistema distribuido, así que como decisión de grupo se decidió utilizar el sistema de Cloud de Google : *Google Cloud Platform*.

El repositorio del proyecto del Motor de Búsqueda se encuentro en GitHub en el enlace ¹

Para realizar este proyecto se tomaron los siguientes puntos:

*e-mail:patrick.marquez@ucsp.edu.pe

†e-mail:jean.cornejo@ucsp.edu.pe

¹<https://github.com/patrick03524/Search-Engine-Inverted-Index-Page-Rank-Cloud-Computing>.

1.1 Configuración del Cluster dentro de *Google Cloud Platform*

Para este proyecto se logró una correcta configuración de un cluster utilizando *Google Cloud Platform*, se utilizó principalmente este tutorial: ²

Clústeres

+ CREATE CLUSTER

ACTUALIZAR

▶ INICIAR

■ DETENER

🗑 BORR

Filtro

Busca clústeres y presiona Intro

?

☰

<input type="checkbox"/>	Nombre ↑	Estado	Región	Zona	Total de nodos trabajadores	Eliminar
<input type="checkbox"/>	buscador	<div>Detenido</div>	us-central1	us-central1-c	2	Desa

Figure 1: Vista principal del cluster configurado.

Lo que se logró con esto es la creación de un cluster distribuido en *Google Cloud Platform* como se puede apreciar en la Figura 1 de tal manera que dentro de este cluster exista un nodo maestro y 2 nodos esclavos de la siguiente manera:

SUPERVISIÓN

TRABAJOS

INSTANCIAS DE VM

CONFIGURACIÓN

INTERFACES

Filtro

Filtrar instancias

<div><input checked="" type="radio"/></div>	Nombre <div>↑</div>	Rol
<div><input checked="" type="radio"/></div>	buscador-m	Instancia principal
<div><input checked="" type="radio"/></div>	buscador-w-0	Trabajador
<div><input checked="" type="radio"/></div>	buscador-w-1	Trabajador

Figure 2: Vista de los Nodos del cluster.

Posteriormente se realizó una configuración de las instancias de las máquinas virtuales de todos los nodos para la configuración del protocolo de comunicación SSH, configuración de las variables de entorno para Java e instalación de todas las librerías de Apache Hadoop para la compilación del proyecto del Motor de Búsqueda.

²<http://www-scf.usc.edu/~shin630/Youngmin/files/HadoopInvertedIndexV5.pdf>.

Por último luego de la configuración y actualización de los sub-sistemas de Linux de nuestras instancias de las Máquinas virtuales, se realizó un lanzamiento de un código de prueba para comprobar que todo funcionará; así que se lanzó una tarea de compilación de un archivo .java dentro del nodo maestro, esto resultó en un archivo .jar así que esto fue enviado con el comando "hadoop fs" para enviarlo a un bucket que nos proporciona *Google Cloud Platform* con el nombre de "dataproc" para cargar nuestra data de manera local. Esto se realizó de igual manera para la compilación del proyecto de Motor de Búsqueda.

1.2 Implementación del Motor de Búsqueda

Dentro de la implementación del Motor de Búsqueda se crearon archivos en Java de clases con el algoritmo de Índice Invertido y el Buscador como tal. Cada clase tenía sus funciones Map y Reduce para lograr la utilización de MapReduce dentro del entorno distribuido en nuestro Clúster. Debido a ciertos problemas no se logró una correcta implementación de los códigos en Java pero los intentos de este proyecto se encuentran en el repositorio adjunto de Github.

Adicionalmente, en el trabajo era necesario añadir como complemento una sección visual de los resultados de nuestros algoritmos, a manera de mostrar un resultado, se ha de resaltar que debido a ciertos problemas con el código y la falta de fiabilidad en los resultados por la incorrecta implementación de las secciones correspondientes al output, no se progresó en la idea original, la cual era usar un framework de desarrollo web y un medio de comunicación para obtener los resultados.

En la imagen 3 se observa los códigos presentes en el sistema de archivos de Google y preparados para ser enviados a Hadoop en la forma de jobs.

Filtrar solo por prefijo de nombre ▼						Filtro Filtrar objetos y carpetas		Mostrar datos borrados	
<input type="checkbox"/>	Nombre	Tamaño	Tipo	Fecha de creación ?	Clase				
<input type="checkbox"/>	Searcher.jar	5.3 KB	application/octet-stream	3 may 2022 23:50:11	Stand...				
<input type="checkbox"/>	Test1.jar	52.2 KB	application/octet-stream	4 may 2022 00:00:02	Stand...				
<input type="checkbox"/>	google-cloud-dataproc-metainfo/	—	Carpeta	—	—				

Figure 3: Códigos incluidos en el clúster preparados para ejecutar jobs.

En la imagen 4 se encuentra el resultado luego de ejecutar alguno de los .jar previamente mostrados dados un tiempo de ejecución.

Google Cloud Platform

My First Project

Buscar Productos, recursos, documentos (f)

Trabajos alojados en clústeres

Trabajos

Flujos de trabajo

Políticas de ajuste de escala...

Sin servidores

Lotus

Servicios públicos

Intercambio de componentes

Metastore

Workbench

Detalles del trabajo

CLONAR BORRAR DETENER ACTUALIZAR

Job failed with message [Exception in thread "main" java.lang.NoClassDefFoundError: Searcher (wrong name: com/mycompany/searcher/Searcher)]. Additional details can be found at: https://console.cloud.google.com/dataproc/jobs/job-507teasf/project-herc...

ID de trabajo

UIDO de trabajo

Tipo

Estado

SUPERVISIÓN

CONFIGURACIÓN

En los siguientes gráficos, se representan las métricas del clúster en el que se ejecutó este trabajo, limitadas al momento en que se ejecutó este trabajo. Es posible que más de un trabajo se ejecute en un clúster a la vez, por lo que es posible que estas métricas no reflejen con precisión el uso de los recursos de este trabajo. Las métricas de un trabajo pueden retrasarse varios minutos en comparación con la ejecución del trabajo.

RESTABLECER EL ZOOM

1 hora 6 horas 12 horas 1 día 2 días 4 días 7 días 14 días 30 días 23:46 - 23:01

Resultado

UNIÓN DE LÍNEA DESACTIVADA

```
Exception in thread "main" java.lang.NoClassDefFoundError: Searcher (wrong name: com/mycompany/searcher/Searcher)
    at java.lang.ClassLoader.defineClass(Native Method)
    at java.lang.ClassLoader.defineClass(ClassLoader.java:756)
    at java.security.SecureClassLoader.defineClass(SecureClassLoader.java:142)
    at java.net.URLClassLoader.defineClass(URLClassLoader.java:473)
    at java.net.URLClassLoader.access$800(URLClassLoader.java:74)
    at java.net.URLClassLoader$1.run(URLClassLoader.java:358)
    at java.net.URLClassLoader$1.run(URLClassLoader.java:353)
    at java.security.AccessController.doPrivileged(Native Method)
    at java.net.URLClassLoader.findClass(URLClassLoader.java:362)
    at java.lang.ClassLoader.loadClass(ClassLoader.java:418)
    at sun.misc.Launcher$BootstrapClassLoader.loadClass(Launcher.java:352)
    at java.lang.ClassLoader.loadClass(ClassLoader.java:355)
    at java.lang.Class.forName0(Native Method)
    at java.lang.Class.forName(Class.java:264)
    at com.google.cloud.hadoop.service.agent.job.shim.HadoopRunClassShim.main(HadoopRunClassShim.java:18)
```

Se completó el resultado

El trabajo job-507teasf se envió correctamente.

Operaciones de My First Project cargadas

Completada

Figure 4: Error mostrado a la hora de ejecutar jobs.

4

1.3 Conclusiones

Como principal conclusión de este trabajo podemos afirmar que la utilización de un sistema distribuido para realizar algoritmos con una gran cantidad de datos es la manera más óptima de realizar y gracias a plataformas como la que provee Google, se nos permite una rápida configuración del cluster como tal para realizar estas tareas.

Posteriormente podemos destacar las dificultades de la implementación de los algoritmos ya que este paradigma distribuido es bastante complicado de manejar pero que provee bastantes ventajas para procesar mucha información. De igual manera *Google Cloud Platform* provee un sistema confiable y robusto con el que se pueden crear algoritmos del Índice Invertido y Page Rank de manera mucho más eficiente.