

Reporte técnico y Reporte de Avance, fases

Tópicos en Telemática

Profesor:

Edwin Montoya

por:

Carlos Mario Blanco Pérez

Santiago Vásquez Zuluaga

Eafit

2017

Procesos:

ETL: En el proceso de extracción de los datos se procederá a descargar el *dataset* de la biblioteca de Gubenberg el cual incluye una colección de documentos tanto en inglés como en español. (Dado que el *dataset* se encuentra en la carpeta raíz su avance es del 100%)

En el proceso de transformación se llevarán todos los elementos que integran el *dataset* a un solo formato .txt (dado que los documentos dentro del *dataset* tienen todos una extensión .txt su avance es del 100%)

En el proceso de carga se llevará el *dataset* a HDFS (Hadoop Distributed File System, cuyo entorno ya está configurado por defecto con un master en la dirección 10.131.137.188 y sus respectivos slaves en los cuales estará distribuida la carga), avance 100%

Procesamiento: El procesamiento se realizará bajo una arquitectura de datos en *batch* donde se integrará al lenguaje de programación *Python* junto con una librería *mrjob* (*map/reduce*) para el manejo y procesamiento del *dataset* bajo un *índice invertido* para clasificar, ordenar y consultas posteriores. Avance 100%

Aplicación: Se desarrollará una aplicación en *Ruby on Rails (RoR)*, la cual será la que tendrá interacción directa con el usuario que en forma de búsqueda el agregará en un *text área* una serie de palabras las cuales retornarán como resultado el documento(s) .txt los cuales tengan la mayor ocurrencia de las palabras ingresadas por el usuario. Avance 100%.

Github: <https://github.com/santivasquez/bigDataProyecto3.git>

Algoritmo Map/reduce, programado en Python

```
1  # -*- coding: utf-8 -*-
2  from mrjob.compat import jobconf_from_env
3  from mrjob.job import MRJob
4  from pymongo import MongoClient
5
6  # PASO 1: Conexión al Server de MongoDB Pasandole el host y el puerto
7  mongoClient = MongoClient('10.131.137.188',27017)
8
9  # PASO 2: Conexión a la base de datos
10 db = mongoClient["grupo_07"]
11 db.authenticate("user1", "eafit.2017")
12
13 # PASO 3: Obtenemos una coleccion para trabajar con ella
14 collection = db.gutenberg
15
16 class MRWordFrequencyCount(MRJob):
17
18     def mapper(self, _, line):
19         value = line.decode('utf-8','ignore').split()
20
21         for word in value:
22             yield (jobconf_from_env('mapreduce.map.input.file'),word), 1
23
24     def reducer(self, key, values):
25         collection.insert({'nameFile': key[0], 'word': key[1], 'numero':sum(values)})
26         yield key,sum(values)
27
28 if __name__ == '__main__':
29     MRWordFrequencyCount.run()
30
31 # PASO FINAL: Cerrar la conexion
32 mongoClient.close()
```

La estructura y outputs en la base de datos Mongoddb son:

```
> db.gutenberg.find({"word": "todos"}).limit(2).pretty()
{
  "_id" : ObjectId("59069beed03b4f1fdafa7d09"),
  "word" : "todos",
  "numero" : 115,
  "nameFile" : "hdfs://master:8020/user/st0263/cblanco/gutenberg-txt-es/10293-8.txt"
}
{
  "_id" : ObjectId("59069c70d03b4f1fdafaa9cb"),
  "word" : "todos",
  "numero" : 115,
  "nameFile" : "hdfs://master:8020/user/st0263/cblanco/gutenberg-txt-es/10293.txt"
}
```

Comandos usados

```
python doc_mapreduce.py -r hadoop hdfs:///user/st0263/cblanco/gutenberg-txt-es/*
```

// Correr el algoritmo Map/reduce pasandole como parametros todos los textos planos dentro de la carpeta Gutenberg-txt-es

```
mongo grupo_07 -u user1 -p
```

// Conectarse a la base de datos MongoDB

```
hadoop fs -ls "hdfs://master:8020/user/st0263/cblanco/"
```

//Inspeccionar

```
hadoop fs -put /home/cblanco/gutenberg-txt-es /user/st0263/cblanco
```

//agregar

```
hadoop fs -cat "hdfs://master:8020/user/st0263/cblanco/gutenberg-txt-es/10293.txt"
```

//Mostrar en pantalla el contenido de "10293.txt"

Configuración – MongoDB – Ruby on Rails – mongoid.yml

```
# (required).
database: grupo_07
# Provides the hosts the default client can connect to.
# of host:port pairs. (required)
hosts:
  - 10.131.137.188:27017

options:
  # Change the default write concern. (default = :majority)
  # write:
  #   w: 1

  # Change the default read preference. Valid options are:
  # :secondary_preferred, :primary, :primary_preferred
  # (default: primary)
  # read:
  #   mode: :secondary_preferred
  #   tag_sets:
  #     - use: web

  # The name of the user for authentication.
  user: 'user1'

  # The password of the user for authentication.
  password: 'eafit.2017'

  # The user's database roles.
```

Indicaciones:

Para iniciar la aplicación, se debe ingresar por la URL 10.131.137.166:3000. Donde se encuentra con una aplicación web elaborada en Ruby on Rails. La interacción se hace mediante el botón de búsqueda, donde se ingresan las palabras a buscar y a continuación se listarán los primeros 20 archivos (traídos de la biblioteca Gutenberg) en donde más se repitan las palabras ingresadas. Los archivos listados contienen un enlace hacia su ubicación, donde se puede tener más información del archivo.

URL: 10.131.137.166:3000