



Documentación Proyecto Final
Product Development
Dashboard Covid-19

Maestría Data Science

Cuarto Trimestre 2020

Número grupo	Nombre Integrante	Carné
4	Elda Magally Calderon Motta	16003182
	Ligia Deyanira Aguilar Estrada	19004789
	Victor Rene Pérez Mayen	19005180

1 Objetivos

- Aplicar conocimientos adquiridos en la clase magistral para utilizar técnicas de diseño e implementación avanzadas en la construcción de un dashboard funcional siguiendo las especificaciones de herramientas y características dadas.
- Combinar el uso de herramientas para el desarrollo de software como lo es el uso de lenguajes de programación (R y Python), bases de datos (PostgreSQL y MySQL) y herramientas de ETL (Airflow).
- Utilizar técnicas de análisis y procesamiento de datos para el correcto despliegue de la información presentada en el dashboard.

2 Resumen

El presente proyecto consiste en poder visualizar la cantidad de casos confirmados, casos fallecidos y casos recuperados del virus COVID-19. Para ello se utilizó un dashboard de Shiny Apps creado en un contenedor de Docker, este dashboard despliega los datos almacenados en una Base de Datos de MySQL la cual tiene comunicación con nuestro dashboard, dicha base de datos tiene la característica que también se encuentra ubicada en un contenedor de Docker; dicha base de datos recibe una ingesta de datos por medio de Airflow , ubicado entro del contenedor de Docker ya mencionado. La librería de Airflow contiene funciones para trabajar y desarrollar procesos de ETL para alimentar la base de datos. Se obtuvo como resultado una visualización a nivel global de los casos de COVID-19, el cual puede recibir información de manera dinámica gracias a las configuraciones y manejo de DAGS (Directed Acyclic Graph por sus siglas en inglés) de Airflow.

3 Introducción

Docker es una aplicación que simplifica el proceso de gestionar los procesos de aplicaciones en contenedores. Los contenedores le permiten ejecutar sus aplicaciones en procesos aislados de recursos. Se parecen a las máquinas virtuales, sin embargo, los contenedores son más portátiles, tienen más recursos y son más dependientes del sistema operativo host.

Para el proyecto Docker toma un papel fundamental en su desarrollo debido a su portabilidad y fácil manejo de contenedores en los cuales se pueden combinar varias herramientas y recursos de software sin que estos tengan conflictos de compatibilidad. Docker adicionalmente es sencillo de mantener y aporta al proyecto servicios de red, manejo optimizado de los recursos a utilizar y autonomía.

A partir de Docker y Docker Compose se manejan diferentes servicios para que nuestro dashboard los utilice y se puedan realizar las diferentes tareas de carga y visualización de datos. Para este proyecto se combinan lenguajes de programación como lo es R y Python, manejo y mantenimiento de base de datos.

En el presente trabajo se explicará cada uno de los servicios contenidos en Docker así como los resultados y conclusiones obtenidos al desarrollar nuestro dashboard.

4 Desarrollo del Dashboard

4.1 Docker y Docker Compose

4.1.1 Dockerfile

Utilizamos un Dockerfile para que pueda crear de forma automática las imágenes necesarias con las instrucciones descritas en el. Se puede decir que es un archivo compuesto por los comandos que el usuario tendría que ingresar manualmente en la terminal de su máquina. Al ejecutar las instrucciones de este archivo se ejecutan de forma secuencial y automatizada.

A continuación se puede observar las instrucciones a ejecutar en nuestro Dockerfile.

```
1 # VERSION 1.10.9
2 # AUTHOR: Matthieu "Puckel_" Roisil
3 # DESCRIPTION: Basic Airflow container
4 # BUILD: docker build --rm -t puckel/docker-airflow .
5 # SOURCE: https://github.com/puckel/docker-airflow
6
7 FROM python:3.7-slim-buster
8 LABEL maintainer="Puckel_"
9
10 # Never prompt the user for choices on installation/configuration
    of packages
11 ENV DEBIAN_FRONTEND noninteractive
12 ENV TERM linux
13
14 # Airflow
15 ARG AIRFLOW_VERSION=1.10.9
16 ARG AIRFLOW_USER_HOME=/usr/local/airflow
17 ARG AIRFLOW_DEPS=""
18 ARG PYTHON_DEPS=""
19 ENV AIRFLOW_HOME=${AIRFLOW_USER_HOME}
20
21 # Define en_US.
22 ENV LANGUAGE en_US.UTF-8
23 ENV LANG en_US.UTF-8
24 ENV LC_ALL en_US.UTF-8
25 ENV LC_CTYPE en_US.UTF-8
26 ENV LC_MESSAGES en_US.UTF-8
27
```

```

28 COPY requirements.txt .
29
30 # Disable noisy "Handling signal" log messages:
31 # ENV GUNICORN_CMD_ARGS --log-level WARNING
32
33 RUN set -ex \
34     && buildDeps=' \
35         freetds-dev \
36         libkrb5-dev \
37         libsasl2-dev \
38         libssl-dev \
39         libffi-dev \
40         libpq-dev \
41         git \
42     ' \
43     && apt-get update -yqq \
44     && apt-get upgrade -yqq \
45     && apt-get install -yqq --no-install-recommends \
46         $buildDeps \
47         freetds-bin \
48         build-essential \
49         default-libmysqlclient-dev \
50         apt-utils \
51         curl \
52         rsync \
53         netcat \
54         locales \
55     && sed -i 's/^# en_US.UTF-8 UTF-8$/en_US.UTF-8 UTF-8/g' /etc/
    locale.gen \
56     && locale-gen \
57     && update-locale LANG=en_US.UTF-8 LC_ALL=en_US.UTF-8 \
58     && useradd -ms /bin/bash -d ${AIRFLOW_USER_HOME} airflow \
59     && pip install -U pip setuptools wheel \
60     && pip install 'SQLAlchemy==1.3.15' \
61     && pip install pytz \
62     && pip install pyOpenSSL \
63     && pip install ndg-httpsclient \
64     && pip install pyasn1 \
65     && pip install apache-airflow[crypto,celery,postgres,hive,jdbc,
    mysql,ssh${AIRFLOW_DEPS:+,}${AIRFLOW_DEPS}]==${AIRFLOW_VERSION}
    \
66     && pip install 'redis==3.2' \
67     && pip install -r requirements.txt \
68     && if [ -n "${PYTHON_DEPS}" ]; then pip install ${PYTHON_DEPS};
    fi \
69     && apt-get purge --auto-remove -yqq $buildDeps \
70     && apt-get autoremove -yqq --purge \
71     && apt-get clean \
72     && rm -rf \
73         /var/lib/apt/lists/* \
74         /tmp/* \
75         /var/tmp/* \
76         /usr/share/man \
77         /usr/share/doc \
78         /usr/share/doc-base
79
80 COPY script/entrypoint.sh /entrypoint.sh

```

```

81 COPY config/airflow.cfg ${AIRFLOW_USER_HOME}/airflow.cfg
82
83 RUN chown -R airflow: ${AIRFLOW_USER_HOME}
84 RUN mkdir -p /home/airflow/monitor
85 RUN chown -R airflow: /home/airflow/monitor
86
87 EXPOSE 8080 5555 8793 3306 33061
88
89 USER airflow
90 WORKDIR ${AIRFLOW_USER_HOME}
91 ENTRYPOINT ["/entrypoint.sh"]
92 CMD ["webserver"]

```

4.1.2 Docker compose

Compose es una herramienta para la cual permite ejecutar aplicaciones en Docker con varios contenedores y servicios corriendo. Con Compose, usa un archivo YAML para configurar los servicios de nuestro proyecto. Permite con un solo comando pueda crear e iniciar todos los servicios con las configuraciones definidas.

El archivo .YAML de nuestro proyecto levanta tres servicios que son:

- Airflow
- Postgres (Base de datos relacionada a Airflow)
- MySQL
- Shiny Server

En cada sección se verá como se encuentran configurados dentro del archivo docker-compose.yaml

4.2 Airflow

4.2.1 Configuración Airflow en docker compose

```

1  postgres:
2    image: postgres:9.6
3    environment:
4      - POSTGRES_USER=airflow
5      - POSTGRES_PASSWORD=airflow
6      - POSTGRES_DB=airflow
7    logging:
8      options:
9        max-size: 10m
10       max-file: "3"
11
12  webserver:
13    build: .
14    restart: always
15    depends_on:

```

```

16     - postgres
17     environment:
18     - LOAD_EX=n
19     - EXECUTOR=Local
20     logging:
21     options:
22         max-size: 10m
23         max-file: "3"
24     volumes:
25     - ./dags:/usr/local/airflow/dags
26     - ./monitor:/home/airflow/monitor
27     # - ./plugins:/usr/local/airflow/plugins
28     ports:
29     - "8080:8080"
30     command: webserver
31     healthcheck:
32     test: [ "CMD-SHELL", "[ -f /usr/local/airflow/airflow-
33         webserver.pid ]" ]
34     interval: 30s
35     timeout: 30s
36     retries: 3

```

4.2.2 DAGS

En Airflow, es un conjunto de tareas que se desean ejecutar, todas estas tareas se encuentran relacionadas unas con otras y se encuentran organizadas para determinar el momento o las acciones que hacen que cada tarea se ejecute.

En esta situación cada DAG se encuentra definido en un archivo de código escrito en el lenguaje de programación Python. A continuación se puede ver cada archivo de Python dentro de la carpeta de DAGS de nuestro proyecto. Cada archivo gobierna y determina las funciones de cada tarea para cargar una tabla diferente para nuestra base de datos.

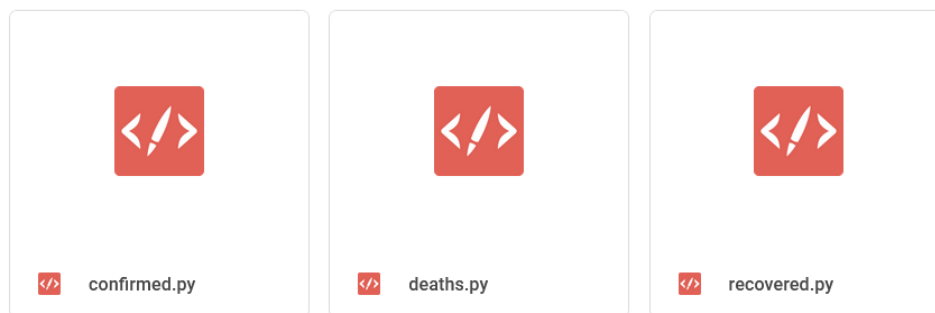


Figure 1: DAGS - Scripts de Python

A continuación se presenta un ejemplo de la construcción de cada DAG para el proyecto.

```

1 import os
2 import pandas as pd
3 from airflow import DAG
4 from airflow.contrib.hooks.fs_hook import FSHook
5 from airflow.contrib.sensors.file_sensor import FileSensor
6 from airflow.hooks.mysql_hook import MySQLHook
7 from airflow.operators.python_operator import PythonOperator
8 from airflow.utils.dates import days_ago
9
10
11 FILE_CONNECTION_ID = "fs_default"
12 FILE_NAME = "time_series_covid19_confirmed_global.csv"
13 OUTPUT_TRANSFORM_FILE = "_confirmed_tmp.csv"
14
15 dag = DAG('confirmed_dag', description='ETL for confirmed table
16         cases',
17         default_args={
18             'owner': 'elda.calderon',
19             'depends_on_past': False,
20             'max_active_runs': 1,
21             'start_date': days_ago(5)
22         },
23         schedule_interval='0 1 * * *',
24         catchup=False)
25
26 file_sensor_task = FileSensor(dag = dag,
27                               task_id = 'file_sensor',
28                               fs_conn_id= FILE_CONNECTION_ID,
29                               filepath= FILE_NAME,
30                               poke_interval= 10,
31                               timeout= 300
32                               )
33
34 def transform_func(**kwargs):
35     folder_path = FSHook(conn_id = FILE_CONNECTION_ID).get_path()
36     file_path = f"{folder_path}/{FILE_NAME}"
37     destination_file = f"{folder_path}/{OUTPUT_TRANSFORM_FILE}"
38     df = pd.read_csv(file_path, header = 0, encoding= 'ISO-8859-1')
39     df.to_csv(destination_file, index = False)
40     os.remove(file_path)
41     return destination_file
42
43 transform_process = PythonOperator(dag = dag,
44                                   task_id = 'transform_process',
45                                   python_callable = transform_func
46                                   ,
47                                   provide_context= True
48                                   )
49
50 def insert_process(**kwargs):
51     ti = kwargs['ti']
52     source_file = ti.xcom_pull(task_ids = 'transform_process')
53     db_connection = MySQLHook('airflow_db').get_sqlalchemy_engine()
54
55     df = pd.read_csv(source_file)

```

```

56     with db_connection.begin() as transaction:
57         transaction.execute("DELETE FROM test.confirmed_table WHERE
58                               1=1")
59         df.to_sql("confirmed_table", con = transaction, schema = "
60                   test", if_exists= "replace", index = False)
61
62     os.remove(source_file)
63
64 insert_process = PythonOperator(dag = dag,
65                                task_id= 'insert_process',
66                                provide_context= True,
67                                python_callable = insert_process
68                                )
69
70 file_sensor_task >> transform_process >> insert_process

```

4.3 Base de Datos - MySQL

MySQL es un sistema de gestión de bases de datos de código abierto. Gestiona sus datos usando una base de datos relacional y SQL (Lenguaje de consulta estructurada).

4.3.1 Configuración MySQL en docker compose

```

1  db:
2      image: mysql:5.7
3      volumes:
4          #- ./db_data:/var/lib/mysql
5          - ./script/schema.sql:/docker-entrypoint-initdb.d/1.sql
6      restart: always
7      ports:
8          - 33061:3306
9      environment:
10         MYSQL_ROOT_PASSWORD: test123
11         MYSQL_DATABASE: test
12         MYSQL_USER: test
13         MYSQL_PASSWORD: test123

```

4.3.2 Tablas y ejecución de scripts

Para la conexión de la base de datos con los contenedores, se procedió a utilizar DataGrip. Se configuró la conexión para MySQL, con valores similares a los que se muestran a continuación:

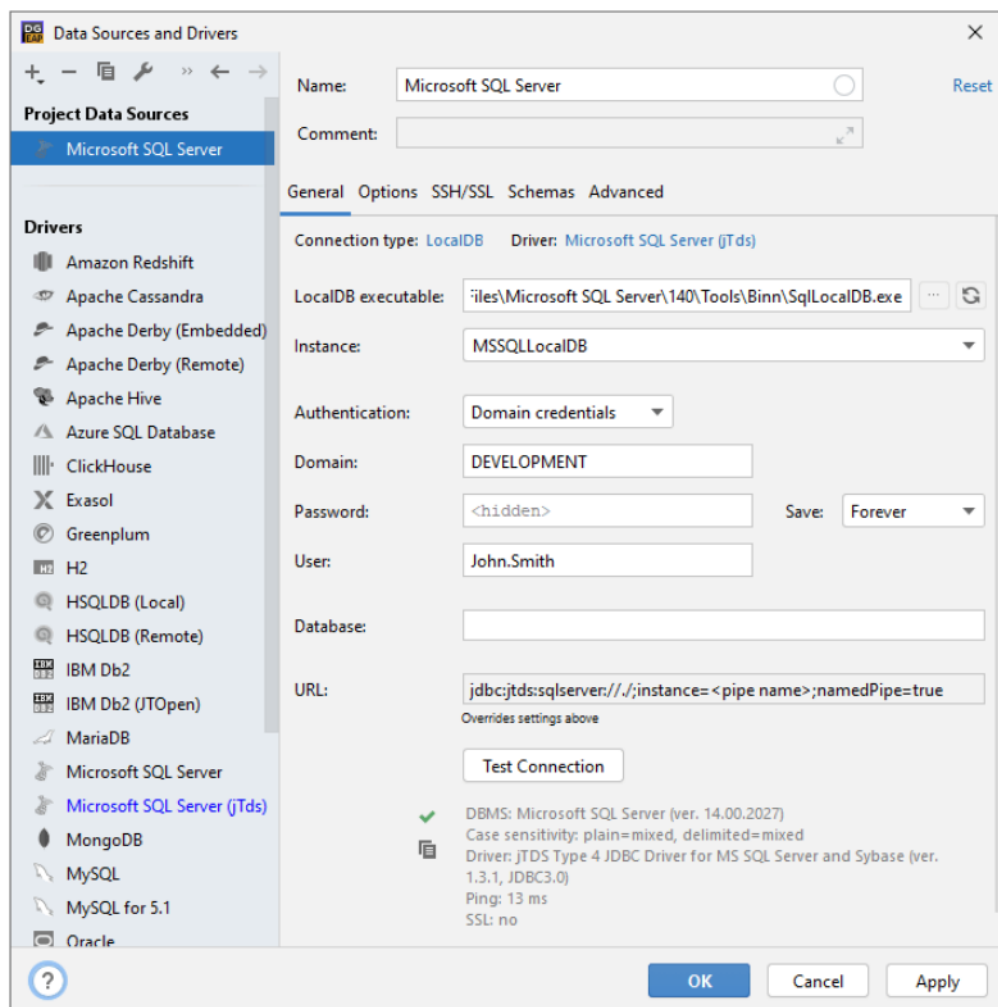


Figure 2: Conexión a la base de datos

Para poder acceder a la base de datos se utilizó MySQL Workbench utilizando los parámetros de conexión que se indicaron al momento de crearla.

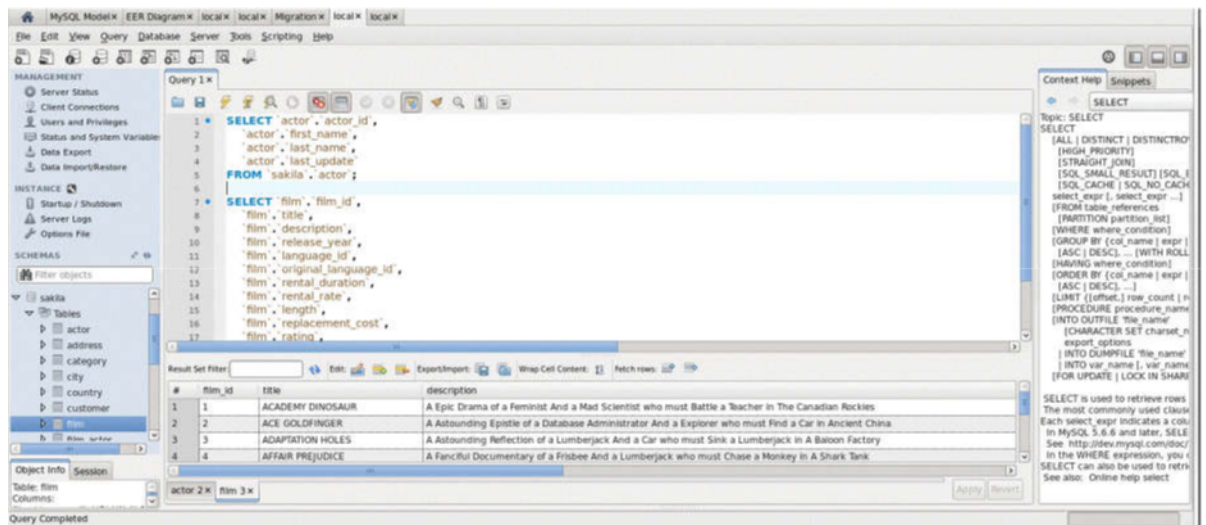


Figure 3: Workbench

4.3.3 Ingesta de información para la base de datos

Por medio de los DAGS de Airflow, se realiza el proceso de ETL hacia la base de datos. Gracias a este proceso, la información de los casos confirmados, recuperados y fallecidos se cargan a la base de datos.

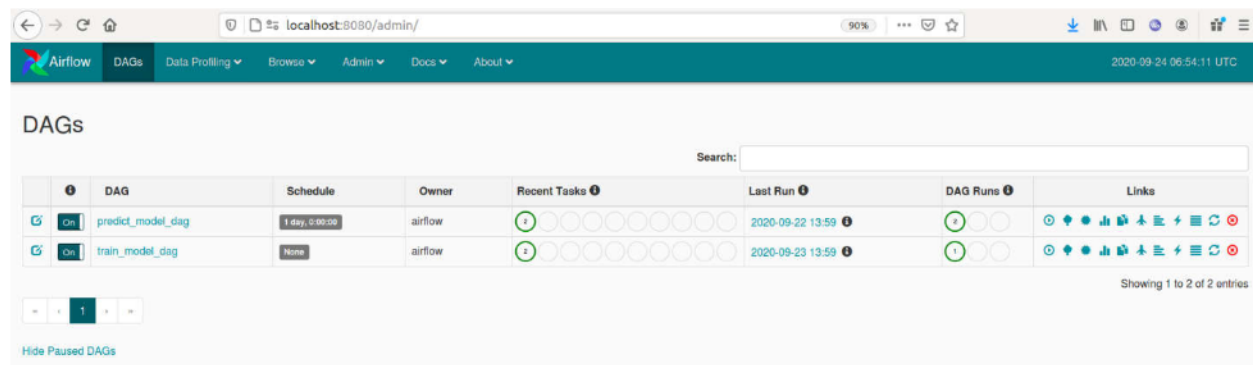


Figure 4: ETL

4.4 Shiny Server

En el apartado de Rsecreo undashboard por medio de Shiny Apps, este dashboard visualiza la cantidad de casos mencionados anteriormente por País, región y por fecha. En el dashboard también se pueden visualizar la cantidad de casos por medio de una gráfica acumulativa, y los 5 países con más casos por fechas específicas.

4.4.1 Configuración Shiny Server en docker compose

```
1 dashboard:
2   build: ./dashboard
3   ports:
4     - 3838:3838 # izquierda puerto maquina, derecha puerto
5   contenedor
6   volumes:
7     - ./prueba:/srv/shiny-server/
8     - ./logs:/var/log/shiny-server/
```

4.4.2 Librerías

```
1 library(shiny)
2 library(leaflet)
3 library(RColorBrewer)
4 library(dplyr)
5 library(RMySQL)
```

4.4.3 Características del Dashboard

Panel con opciones donde el usuario puede determinar el tipo de datos a desplegar.

Options Panel

Select a cathegory:

☐ Confirmed

☐ Deaths

☐ Recovered

Figure 5: Option Panel

Selección de país y provincia (si es que está disponible la información) para el despliegue de las estadísticas y visualización del mapa.

Options Panel

Select a category:

☒ Confirmed
☐ Deaths
☐ Recovered

Select a country

Afghanistan ▼

Select a province

Province.State ▼

Set a Date:

2020-11-23

Number of Cases

44,988

Figure 6: País, provincia y número de casos según las opciones seleccionadas

Gráfica acumulativa de los casos según el país y provincia seleccionados.

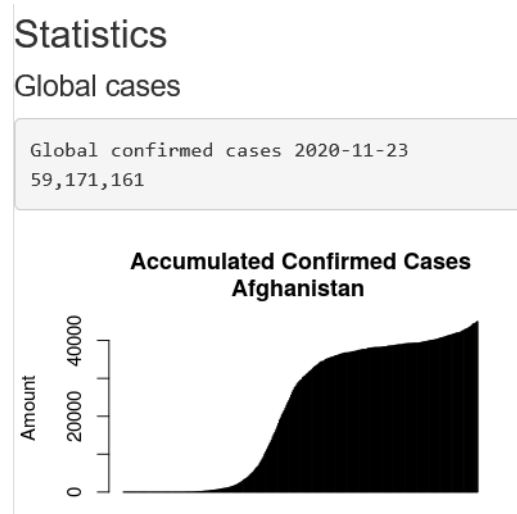


Figure 7: Estadísticas

Top 5 de los países con más casos según la fecha especificada.

Top 5 Country Accumulated Cases

Country.Region	X11.23.20
US	12418311
India	9177840
Brazil	6087608
France	2144979
Russia	2096749

Figure 8: Top 5

Mapa interactivo en el cual se puede ver la información por país.

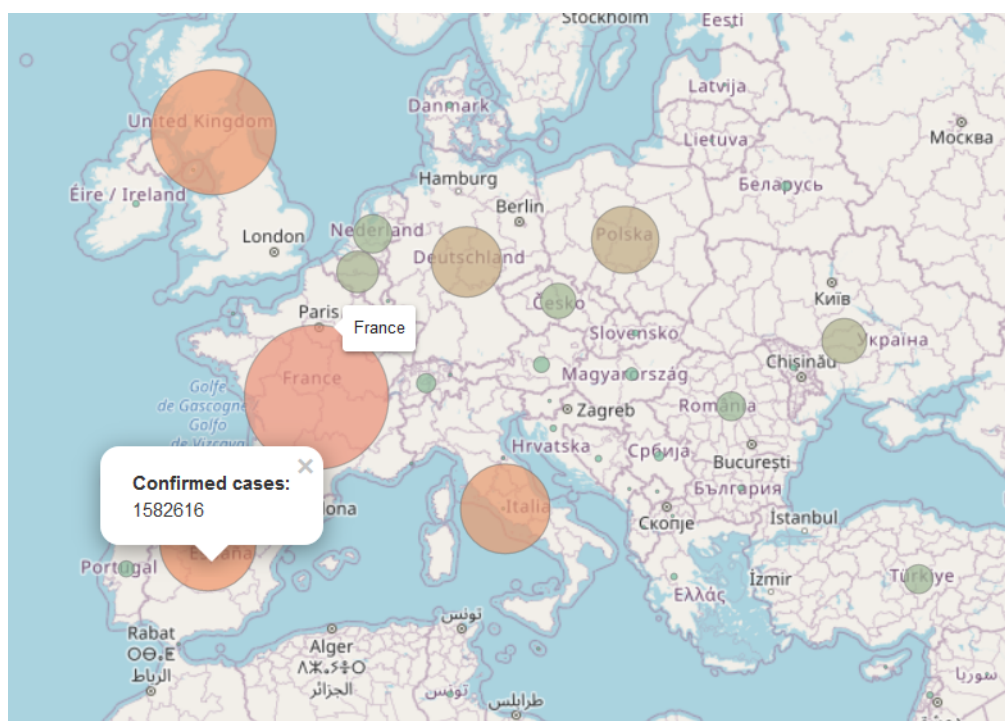


Figure 9: Mapa interactivo

5 Justificación

El proyecto presentado se realizó con el fin de poder realizar una conexión entre los DAGS de Airflow con una base de datos con contenedores hechos en Docker, y que los datos se puedan visualizar de una manera ordenada y simplificada para que al usuario final pueda sacar sus propias conclusiones o inferencias a partir de la data mostrada. Utilizando dashboards elaborados en Shiny Apps de R en conjunto con las otras herramientas descritas anteriormente el fin principal es que toda la aplicación sea portátil, fácil de mantener y el usuario pueda visualizar de manera asequible por medio de una virtualización óptima en cuanto al uso de recursos gracias a los contenedores administrados en Docker.

El despliegue de los casos recuperados, confirmados y fallecidos a causa del virus de COVID-19 son visualizados de tal forma que se espera que contribuya al entendimiento de este problema que nos ha afectado durante el presente año.

6 Recomendaciones

- Se recomienda utilizar contenedores de Docker para base de datos MySQL o SQL Server, ya que son los que tienen mayor cantidad de documentación.
- Es recomendable corroborar la correcta conexión entre los dashboards con los contenedores, para que la data sea visualizada en tiempo real.
- Se debe tener especial cuidado en la configuración de los DAGS para que no tenga errores el proceso de ETL en los contenedores.

7 Conclusiones

- Se instaló Docker, se trabajó con imágenes y contenedores haciendo push del material visto en el curso.
- Ahora se tiene una configuración de MySQL para base de datos la cual fácilmente puede ser adaptada para el desarrollo de otros proyectos.
- Se realizó una conexión de la base de datos con los contenedores de Docker usando DataGrip
- Se hizo la conexión de R server con la base de datos, ambos levantados en Docker container.
- Se demostró que Airflow es permite el proceso de ETL hacia la base de datos.