# COVID19 Data Analysis from Google Open Data using Hive and Tableau

**Authors:** Suleima Gonzalez, Bryan Chica, Mohammed Shodimu, Phillip Navarro
**Instructor:** Jongwook Woo

## Objectives
- Upload data and analyze via wget
- Create a table and test data via Hive
- Clean and test data via Pig
- Store and Export data into Hadoop then Local Machine
- Visualizations and Analysis using PowerBI

## Platform Specifications
- **Hostname:** bigdaiun0.sub03291929060.trainingvcn.oraclevcn.com
- **OS:** Oracle Linux Server 7.9
- **Storage:** 120 GB
- **Memory Size:** 32 GB
- **CPU:** AMC EPYX 7J13 64-Core Processor @ 2.45 Ghz
- **Cluster Version:** Hadoop 3.1.2
- **Cluster Number of Nodes:** 3

**Prerequisites**
- **Operating Systems**
  - Windows 10 or 11
  - MacOS
  - Linux
- **Shell terminal**
  - For windows, download and install Git Bash: *https://git-scm.com/downloads*
  - On Linux or Mac computer, use **terminal**
- **Microsoft Excel**
- **Microsoft PowerBI**
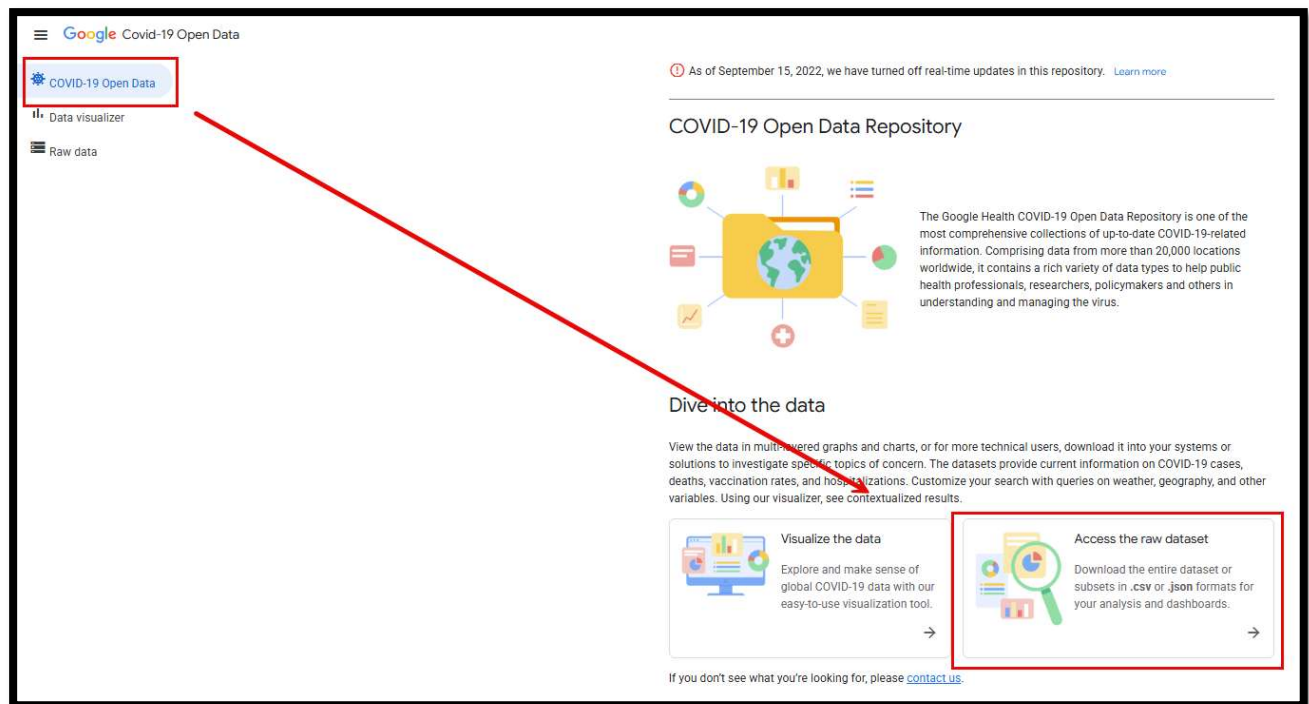- **OpenRefine**

# Introduction

This project demonstrates a complete big data analytic workflow using Google COVID-19 Open Data. The process includes data acquisition, data cleaning, distributed storage, table creation using Hive, aggregation of large datasets, and visualization of key results. By following this step-by-step tutorial, we apply core concepts from CIS 4560 related to big data processing, data warehousing, and analytical reporting.

# Step 1: Download datasets from Google Covid-19 Open Data

*The purpose of this step is to acquire raw COVID-19 data from a reliable and publicly available source. Google COVID-19 Open Data provides standardized datasets that include epidemiological, demographic, geographic, and vaccination information across multiple regions. Downloading these datasets serves as the foundation for all subsequent data processing, analysis, and visualization tasks in this project.*

Go to https://health.google.com/covid-19/open-data/



Download the csv files of the following tables:

- Index
- Demographics
- Epidemiology
- Geography
- Health

- Hospitalizations
- Vaccinations
- By age
- By sex



Clean tables using OpenRefine

Selected the fields I want to drop



Export the file into a csv



Do this with all the csv files

# Step 2: Load Data in Oracle Linux Server

*This step transfers the cleaned datasets from the local machine to the Oracle Linux Server environment. Uploading the files to the cluster allows the data to be processed using distributed computing tools such as Hadoop and Hive. Secure Shell (SSH) is used to establish a remote connection to the server and manage files directly within the big data environment.*

Open **Git Bash** terminal on your local computer



In the shell terminal, use ssh to connect to the Oracle cluster

```
ssh username@ip
```

Download the following file "covid19_tables.zip" into oracle big data server

```
wget https://github.com/chica-94/COVID19-Data-Analysis/raw/refs/heads/main/covid19_tables.zip

ls -lha covid19_tables.zip
```



Unzip tables folder

```
unzip covid19_tables.zip

ls -lha covid19_tables
```

# Step 3: Upload datasets to HDFS

*Hadoop Distributed File System (HDFS) is used to store large datasets across multiple nodes for efficient access and fault tolerance. In this step, directories are created for each dataset, and the CSV files are uploaded into HDFS. Organizing the data into separate directories ensures structured storage and prepares the datasets for table creation in Hive.*

Create a directory rating to put the files to HDFS. Then upload the covid data to the new directory in HDFS. Once the file is added, verify *all* files are uploaded to **covid19** directory (table):

```
cd covid19_tables

hdfs dfs -mkdir tmp/group3_covid19/index
hdfs dfs -mkdir tmp/group3_covid19/hospitalizations
hdfs dfs -mkdir tmp/group3_covid19/health
hdfs dfs -mkdir tmp/group3_covid19/vaccinations
hdfs dfs -mkdir tmp/group3_covid19/demographics
hdfs dfs -mkdir tmp/group3_covid19/geography
hdfs dfs -mkdir tmp/group3_covid19/gender
hdfs dfs -mkdir tmp/group3_covid19/epidemiology
hdfs dfs -mkdir tmp/group3_covid19/age
hdfs dfs -mkdir tmp/group3_covid19/gender
hdfs dfs -mkdir /tmp/group3_covid19/agggregated
hdfs dfs -mkdir /tmp/group3_covid19/top10_countries/

hdfs dfs -ls tmp/group3_covid19/
```

Upload the covid19 datasets into their respective directories in HDFS.

```
hdfs dfs -put index.csv tmp/group3_covid19/index
hdfs dfs -put hospitalizations.csv tmp/group3_covid19/hospitalizations
hdfs dfs -put health.csv tmp/group3_covid19/health
hdfs dfs -put vaccinations.csv tmp/group3_covid19/vaccinations
hdfs dfs -put demographics.csv tmp/group3_covid19/demographics
hdfs dfs -put geography.csv tmp/group3_covid19/geography
hdfs dfs -put gender.csv tmp/group3_covid19/gender
hdfs dfs -put epidemiology.csv tmp/group3_covid19/epidemiology
hdfs dfs -put age.csv tmp/group3_covid19/age
```

Ensure that all files are loaded to hdfs and are in their designated directories

```
hdfs dfs -ls -R tmp/group3_covid19/
```

# Step 4: Create Tables Using Beeline

*The purpose of this step is to create external Hive tables that map directly to the datasets stored in HDFS. Using Beeline and HiveQL commands, each dataset is structured into a table with defined schemas. External tables are used so that the underlying data remains in HDFS while still allowing SQL-based querying and analysis.*

Now that the data is loaded in HDFS, we can now use Hive to create a table and query the data using hiveql commands

Open the Hive shell environment using **beeline**.

```
beeline
```

Use your database to create tables

```
use username;
```

Paste the following HiveQL code to create an external table called which will contain all the data from each dataset uploaded to HDFS. We start by applying a "DROP" command to remove any existing table that may come into conflict

### Create **covid19_index** table

```
DROP TABLE IF EXISTS covid19_index;

CREATE EXTERNAL TABLE IF NOT EXISTS covid19_index(
        location_key STRING,
        country_code STRING,
        country_name STRING,
        subregion1_name STRING,
        subregion2_name STRING)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
STORED AS TEXTFILE LOCATION '/tmp/group3_covid19/index/'
TBLPROPERTIES ("skip.header.line.count"="1");
```

### Create **covid19_hospitalizations** table

```
DROP TABLE IF EXISTS covid19_hospitalizations;

CREATE EXTERNAL TABLE IF NOT EXISTS covid19_hospitalizations(
        event_date DATE,
        location_key STRING,
        new_hospitalized_patients INT,
        cumulative_hospitalized_patients INT,
        new_intensive_care_patients INT,
        cumulative_intensive_care_patients INT)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE LOCATION '/tmp/group3_covid19/hospitalizations/'
TBLPROPERTIES ("skip.header.line.count"="1");
```

## Create **covid19_health** table

```
DROP TABLE IF EXISTS covid19_health;

CREATE EXTERNAL TABLE IF NOT EXISTS covid19_health(
      location_key STRING,
      adult_male_mortality_rate DOUBLE,
      adult_female_mortality_rate DOUBLE,
      life_expectancy DOUBLE,
      diabetes_prevalence DOUBLE,
      health_expenditure_usd DOUBLE,
      out_of_pocket_health_expenditure_usd DOUBLE)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
STORED AS TEXTFILE LOCATION '/tmp/group3_covid19/health/'
TBLPROPERTIES ("skip.header.line.count"="1");
```

## Create **covid19_vaccinations** table

```
DROP TABLE IF EXISTS covid19_vaccinations;

CREATE EXTERNAL TABLE IF NOT EXISTS covid19_vaccinations(
      event_date DATE,
      location_key STRING,
      new_persons_fully_vaccinated INT,
cumulative_persons_fully_vaccinated INT)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
STORED AS TEXTFILE LOCATION '/tmp/group3_covid19/vaccinations/'
TBLPROPERTIES ("skip.header.line.count"="1");
```

## Create **covid19_geography** table

```
DROP TABLE IF EXISTS covid19_geography;

CREATE EXTERNAL TABLE IF NOT EXISTS covid19_geography(
      location_key STRING,
      latitude DOUBLE,
      longitude DOUBLE)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
STORED AS TEXTFILE LOCATION '/tmp/group3_covid19/geography/'
TBLPROPERTIES ("skip.header.line.count"="1");
```

## Create **covid19_demographics** table

```
DROP TABLE IF EXISTS covid19_demographics;

CREATE EXTERNAL TABLE IF NOT EXISTS covid19_demographics(
       location_key STRING,
       population INT,
       population_male INT,
       population_female INT,
       population_age_00_09 INT,
       population_age_10_19 INT,
       population_age_20_29 INT,
       population_age_30_39 INT,
       population_age_40_49 INT,
       population_age_50_59 INT,
       population_age_60_69 INT,
       population_age_70_79 INT,
       population_age_80_and_older INT)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
STORED AS TEXTFILE LOCATION '/tmp/group3_covid19/demographics/'
TBLPROPERTIES ("skip.header.line.count"="1");
```

## Create **covid19_gender** table

```
DROP TABLE IF EXISTS covid19_gender;

CREATE EXTERNAL TABLE IF NOT EXISTS covid19_gender(
       event_date DATE,
       location_key STRING,
       new_deceased_male INT,
       cumulative_deceased_male INT,
       new_deceased_female INT,
       cumulative_deceased_female INT,
       new_recovered_male INT,
       cumulative_recovered_male INT,
       new_recovered_female INT,
       cumulative_recovered_female INT)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
STORED AS TEXTFILE LOCATION '/tmp/group3_covid19/gender/'
TBLPROPERTIES ("skip.header.line.count"="1");
```

## Create **covid19_epidemiology** table

```
DROP TABLE IF EXISTS covid19_epidemiology;

CREATE EXTERNAL TABLE IF NOT EXISTS covid19_epidemiology(
       event_date DATE,
       location_key STRING,
       new_confirmed INT,
       cumulative_confirmed INT,
       new_deceased INT,
       cumulative_deceased INT,
       new_recovered INT,
       cumulative_recovered INT)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
STORED AS TEXTFILE LOCATION '/tmp/group3_covid19/epidemiology/'
TBLPROPERTIES ("skip.header.line.count"="1");
```

Create **covid19_age** table

```
DROP TABLE IF EXISTS covid19_age;

CREATE EXTERNAL TABLE IF NOT EXISTS covid19_age(
       event_date DATE,
       location_key STRING,
       new_recovered_age_0-9 INT,
       new_recovered_age_10-19 INT,
       new_recovered_age_20-29 INT,
       new_recovered_age_30-39 INT,
       new_recovered_age_40-49 INT,
       new_recovered_age_50-59 INT,
       new_recovered_age_60-69 INT,
       new_recovered_age_70-79 INT,
       new_recovered_age_80-89 INT,
       new_recovered_age_90-99 INT,
       cumulative_recovered_age_0-9 INT,
       cumulative_recovered_age_10-19 INT,
       cumulative_recovered_age_20-29 INT,
       cumulative_recovered_age_30-39 INT,
       cumulative_recovered_age_40-49 INT,
       cumulative_recovered_age_50-59 INT,
       cumulative_recovered_age_60-69 INT,
       cumulative_recovered_age_70-79 INT,
       cumulative_recovered_age_80-89 INT,
       cumulative_recovered_age_90-99 INT )
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
STORED AS TEXTFILE LOCATION '/tmp/group3 covid19/age/'
```

Then, in the Hive shell, make sure that the tables **covid19** are created and is shown:

# Step 4: Create external aggregate table

*An aggregated table is created to combine multiple COVID-19 datasets into a single unified structure. This table enables comprehensive analysis by integrating epidemiology, hospitalizations, vaccinations, health indicators, demographics, and geographic data. Aggregating the data simplifies querying and improves analytical efficiency.*

Paste the following HiveQL code to create a table called "**covid19_aggregated**," to store data from the **covid19** tables.

```
DROP TABLE IF EXISTS covid19_aggregated;

CREATE EXTERNAL TABLE IF NOT EXISTS covid19_aggregated(
        event_date DATE, location_key STRING,country_code STRING,country_name
        STRING, subregion1_name STRING,subregion2_name STRING, new_confirmed INT,
        cumulative_confirmed INT,new_deceased INT,cumulative_deceased
        INT,new_recovered INT,cumulative_recovered INT,new_hospitalized_patients
        INT,cumulative_hospitalized_patients INT,new_intensive_care_patients
        INT,cumulative_intensive_care_patients INT,new_persons_fully_vaccinated
        INT,cumulative_persons_fully_vaccinated INT,adult_male_mortality_rate
        DOUBLE,adult_female_mortality_rate DOUBLE,life_expectancy
        DOUBLE,diabetes_prevalence DOUBLE,health_expenditure_usd
        DOUBLE,out_of_pocket_health_expenditure_usd DOUBLE,population
        INT,population_male INT,population_female INT,population_age_00_09
        INT,population_age_10_19 INT,population_age_20_29 INT,population_age_30_39
        INT,population_age_40_49 INT,population_age_50_59 INT,population_age_60_69
        INT,population_age_70_79 INT,population_age_80_and_older INTlatitude
        DOUBLE,longitude DOUBLE,new_deceased_male INT,cumulative_deceased_male
        INT,new_deceased_female INT,cumulative_deceased_female
        INT,new_recovered_male INT,cumulative_recovered_male
        INT,new_recovered_female INT,cumulative_recovered_female
        INTnew_recovered_age_0-9 INT,new_recovered_age_10-19
        INT,new_recovered_age_20-29 INT,new_recovered_age_30-39
        INT,new_recovered_age_40-49 INT,new_recovered_age_50-59
        INT,new_recovered_age_60-69 INT,new_recovered_age_70-79
        INT,new_recovered_age_80-89 INT,new_recovered_age_90-99
        INT,cumulative_recovered_age_0-9 INT,cumulative_recovered_age_10-19
        INT,cumulative_recovered_age_20-29 INT,cumulative_recovered_age_30-39
        INT,cumulative_recovered_age_40-49 INT,cumulative_recovered_age_50-59
        INT,cumulative_recovered_age_60-69 INT,cumulative_recovered_age_70-79
        INT,cumulative_recovered_age_80-89 INT,cumulative_recovered_age_90-99 INT)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
STORED AS TEXTFILE LOCATION '/tmp/group3_covid19/aggregated';
```

Now you need to populate data from the **covid19 tables** to **covid19_aggregated** using SELECT statement as follows:

```
INSERT OVERWRITE TABLE covid19_aggregated
SELECT
        epi.event_date,epi.location_key,idx.country_code,idx.country_name,idx.subregion1_nam
        e,idx.subregion2_name,epi.new_confirmed,epi.cumulative_confirmed,epi.new_deceased,
        epi.cumulative_deceased,epi.new_recovered,epi.cumulative_recovered,hosp.new_hospital
        ized_patients,hosp.cumulative_hospitalized_patients,hosp.new_intensive_care_patients
        ,hosp.cumulative_intensive_care_patients,vac.new_persons_fully_vaccinated,vac.cumula
        tive_persons_fully_vaccinated,health.adult_male_mortality_rate,health.adult_female_m
        ortality_rate,health.life_expectancy,health.diabetes_prevalence,health.health_expend
        iture_usd,health.out_of_pocket_health_expenditure_usd,demo.population,demo.populatio
        n_male,demo.population_female,demo.population_age_00_09,demo.population_age_10_19,de
        mo.population_age_20_29,demo.population_age_30_39,demo.population_age_40_49,demo.pop
        ulation_age_50_59,demo.population_age_60_69,demo.population_age_70_79,demo.populatio
        n_age_80_and_older,geo.latitude,geo.longitude,gender.new_deceased_male,gender.cumula
        tive_deceased_male,gender.new_deceased_female,gender.cumulative_deceased_female,gend
        er.new_recovered_male,gender.cumulative_recovered_male,gender.new_recovered_female,g
        ender.cumulative_recovered_female,age.new_recovered_age_0_9,age.new_recovered_age_10
        _19,age.new_recovered_age_20_29,age.new_recovered_age_30_39,age.new_recovered_age_40
        _49,age.new_recovered_age_50_59,age.new_recovered_age_60_69,age.new_recovered_age_70
        _79,age.new_recovered_age_80_89,age.new_recovered_age_90_99,age.cumulative_recovered
        _age_0_9,age.cumulative_recovered_age_10_19,age.cumulative_recovered_age_20_29,age.c
        umulative_recovered_age_30_39,age.cumulative_recovered_age_40_49,age.cumulative_reco
        vered_age_50_59,age.cumulative_recovered_age_60_69,age.cumulative_recovered_age_70_7
        9,age.cumulative_recovered_age_80_89,age.cumulative_recovered_age_90_99
        FROM covid19_epidemiology epi
LEFT JOIN covid19_index idx ON epi.location_key = idx.location_key
LEFT JOIN covid19_hospitalizations hosp ON epi.location_key = hosp.location_key AND
epi.event_date = hosp.event_date
LEFT JOIN covid19_vaccinations vac ON epi.location_key = vac.location_key AND
epi.event_date = vac.event_date
LEFT JOIN covid19_health health ON epi.location_key = health.location_key
LEFT JOIN covid19_demographics demo ON epi.location_key = demo.location_key
LEFT JOIN covid19_geography geo ON epi.location_key = geo.location_key
LEFT JOIN covid19_gender gender ON epi.location_key = gender.location_key AND
epi.event_date = gender.event_date
LEFT JOIN covid19_age age ON epi.location_key = age.location_key AND epi.event_date =
age.event_date;
```

Check if the table **covid19_aggregated** tables are shown

```
show tables;
```

Query the content of the **covid19_aggregated** table to ensure that the table has been created successfully.

```
SELECT * FROM covid19_aggregated LIMIT 3;
```

# Step 5: Data Cleaning and Preparation

*This step focuses on refining the dataset to support meaningful analysis. By querying the aggregated table, the top ten countries with the highest cumulative confirmed COVID-19 cases are identified. Reducing the dataset to the most relevant records improves performance and allows the analysis to focus on the most impacted regions.*

Query the aggregated tables to find the top 10 countries with the most confirmed COVID-19 cases

```
SELECT country_name, MAX(cumulative_confirmed) AS total_cases
FROM covid19_aggregated
WHERE cumulative_confirmed IS NOT NULL
GROUP BY country_name
ORDER BY total_cases DESC
LIMIT 10;
```

```
+--------------------------+--------------+
|      country_name        | total_cases  |
+--------------------------+--------------+
| United States of America | 92440495     |
| India                    | 44516479     |
| Brazil                   | 34568833     |
| France                   | 33766090     |
| Germany                  | 32604993     |
| South Korea              | 24264470     |
| United Kingdom           | 23554971     |
| Italy                    | 22114423     |
| Russia                   | 20265004     |
| Japan                    | 19868288     |
+--------------------------+--------------+
```

In the Hive shell, use the following code to create an empty table called **covid19_top10_countries**

```
DROP TABLE IF EXISTS covid19_top10_countries;

CREATE EXTERNAL TABLE IF NOT EXISTS covid19_top10_countries (
        location_key STRING, country_code STRING, country_name STRING, subregion1_name
        STRING, subregion2_name STRING, cumulative_confirmed INT, cumulative_deceased INT,
        cumulative_recovered INT, cumulative_hospitalized_patients INT,
        cumulative_intensive_care_patients INT, cumulative_persons_fully_vaccinated INT,
        adult_male_mortality_rate DOUBLE, adult_female_mortality_rate DOUBLE,
        life_expectancy DOUBLE, diabetes_prevalence DOUBLE, health_expenditure_usd DOUBLE,
        out_of_pocket_health_expenditure_usd DOUBLE, population INT, population_male INT,
        population_female INT, population_age_00_09 INT, population_age_10_19 INT,
        population_age_20_29 INT, population_age_30_39 INT,population_age_40_49 INT,
        population_age_50_59 INT, population_age_60_69 INT, population_age_70_79 INT,
        population_age_80_and_older INT, latitude DOUBLE, longitude DOUBLE,
        cumulative_deceased_male INT, cumulative_deceased_female INT,
        cumulative_recovered_male INT, cumulative_recovered_female INT,
        cumulative_recovered_age_0_9 INT, cumulative_recovered_age_10_19 INT,
        cumulative_recovered_age_20_29 INT, cumulative_recovered_age_30_39 INT,
        cumulative_recovered_age_40_49 INT, cumulative_recovered_age_50_59 INT,
        cumulative_recovered_age_60_69 INT, cumulative_recovered_age_70_79 INT,
        cumulative_recovered_age_80_89 INT, cumulative_recovered_age_90_99 INT)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE
LOCATION '/tmp/qroup3 covid19/top10 countries/';
```

Copy and paste the following Hive code to Hive shell in order to create a table
**covid19_top10_countries.** It will be data from the top 10 countries with the most confirmed cases
from the **covid19_aggregated** tables

```
WITH top10 AS (
    SELECT country_name, MAX(cumulative_confirmed) AS total_cases
    FROM covid19_aggregated
    WHERE cumulative_confirmed IS NOT NULL
    GROUP BY country_name
    ORDER BY total_cases DESC
    LIMIT 10)
INSERT OVERWRITE TABLE covid19_top10_countries
SELECT
    a.location_key, a.country_code, a.country_name, a.subregion1_name, a.subregion2_name,
    MAX(a.cumulative_confirmed) AS cumulative_confirmed, MAX(a.cumulative_deceased) AS
    cumulative_deceased, MAX(a.cumulative_recovered) AS cumulative_recovered,
    MAX(a.cumulative_hospitalized_patients) AS cumulative_hospitalized_patients,
    MAX(a.cumulative_intensive_care_patients) AS cumulative_intensive_care_patients,
    MAX(a.cumulative_persons_fully_vaccinated) AS cumulative_persons_fully_vaccinated,
    MAX(a.adult_male_mortality_rate), MAX(a.adult_female_mortality_rate),
    MAX(a.life_expectancy), MAX(a.diabetes_prevalence), MAX(a.health_expenditure_usd),
    MAX(a.out_of_pocket_health_expenditure_usd), MAX(a.population),
    MAX(a.population_male), MAX(a.population_female), MAX(a.population_age_00_09),
    MAX(a.population_age_10_19), MAX(a.population_age_20_29),
    MAX(a.population_age_30_39), MAX(a.population_age_40_49),
    MAX(a.population_age_50_59), MAX(a.population_age_60_69),
    MAX(a.population_age_70_79), MAX(a.population_age_80_and_older), a.latitude,
    a.longitude, MAX(a.cumulative_deceased_male), MAX(a.cumulative_deceased_female),
    MAX(a.cumulative_recovered_male), MAX(a.cumulative_recovered_female),
    MAX(a.cumulative_recovered_age_0_9), MAX(a.cumulative_recovered_age_10_19),
    MAX(a.cumulative_recovered_age_20_29), MAX(a.cumulative_recovered_age_30_39),
    MAX(a.cumulative_recovered_age_40_49), MAX(a.cumulative_recovered_age_50_59),
    MAX(a.cumulative_recovered_age_60_69), MAX(a.cumulative_recovered_age_70_79),
    MAX(a.cumulative_recovered_age_80_89), MAX(a.cumulative_recovered_age_90_99)
FROM covid19_aggregated a
JOIN top10 t ON a.country_name = t.country_name
WHERE a.subregion1_name IS NOT NULL
GROUP BY a.location_key,a.country_code,a.country_name,a.subregion1_name,a.subregion2_name;
```

Query from the table and to see if it has the correct data and values:

```
SELECT * FROM covid19_top10_countries LIMIT 10;
```

In the Beeline shell, you need to check if the table **covid19_top10_countries** shows your
database and hdfs directory

```
describe formatted covid19_top10_countries;
```

Then, in the hive shell, you need to check if the tables **covid19_aggregated** and **covid19_top10_countries** are shown:

```
show tables;
```



Compare the number of records between the tables **covid19_aggregated** and **covid19_top10_countries.** You will see the table was reduced from **12,525,825** to **10,660.**

```
SELECT COUNT(*) FROM covid19_aggregated;
SELECT COUNT(*) FROM covid19_top10_countries;
```



# Step 6: Export Cleaned Table

*After identifying the top ten countries, the cleaned dataset is exported from HDFS to a single CSV file. The (-getmerge) command is used to combine multiple output files into one consolidated file, which is then downloaded to the local machine. This exported dataset is used for visualization and further analysis.*

After the **covid19_top10_countries** table is created, listing the contents inside the **/tmp/group3_covid19/top10_countries/** will show 49 files

```
hdfs dfs -ls /tmp/group3_covid19/top10_countries/
```

Merge the 49 output files to a file named **covid19_top10_countries.csv** using the following:

```
hdfs dfs -getmerge /tmp/group3_covid19/top10_countries/0000*
covid19_top10_countries.csv
```
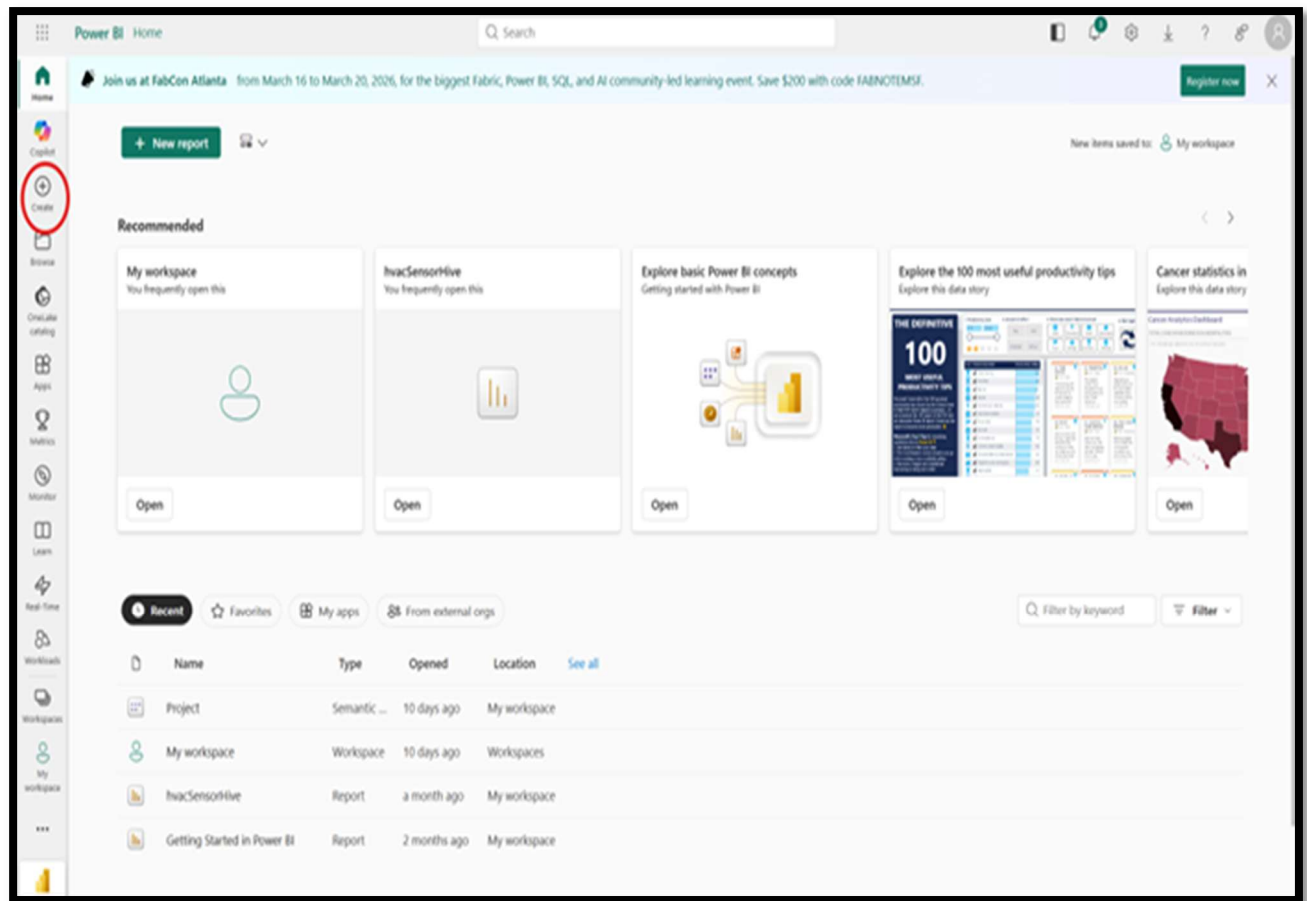
Logout of the cluster and download the csv file into your local machine

```
scp username@ip:/home/bchica/covid19_tables/covid19_top10_countries.csv
~/Downloads/
```
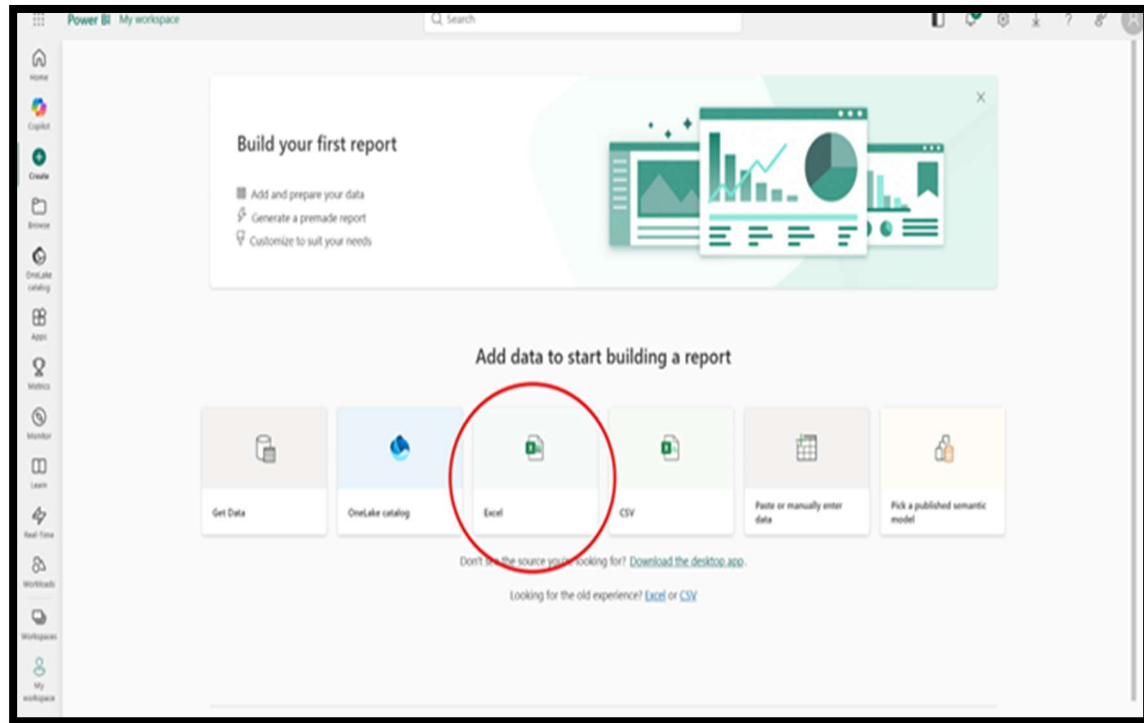
# Step 7: Visualization using PowerBI

*Power BI is used to visualize the results of the aggregated COVID-19 analysis. A bubble map is created to display the top ten countries by cumulative confirmed COVID-19 cases, where bubble size represents the relative magnitude of total confirmed cases. This visualization provides a clear geographic comparison of COVID-19 impact among the most affected countries and supports the findings derived from the Hive-based analysis.*
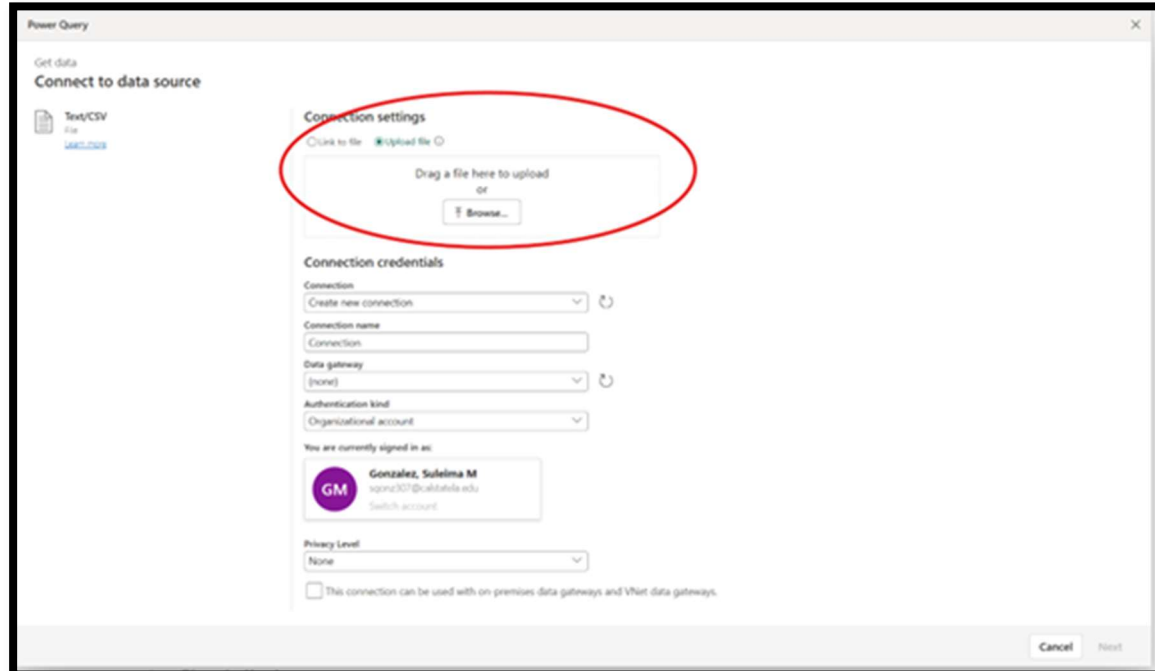
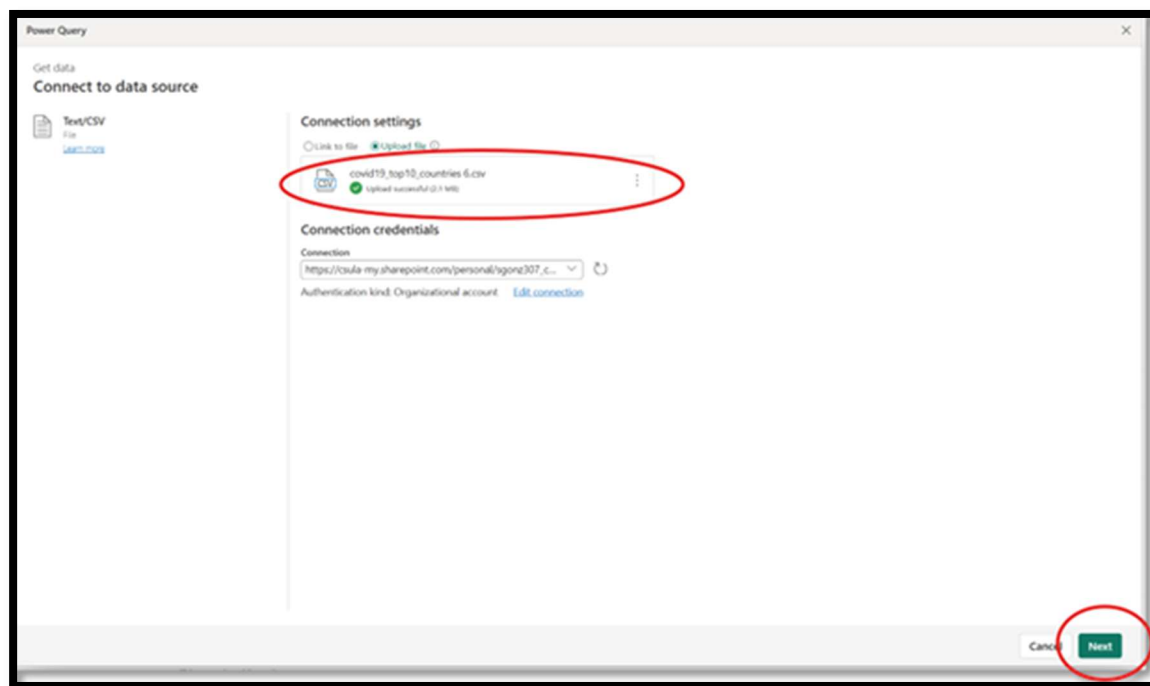Open your MS Power BI Desktop at your local computer. Select **Create**

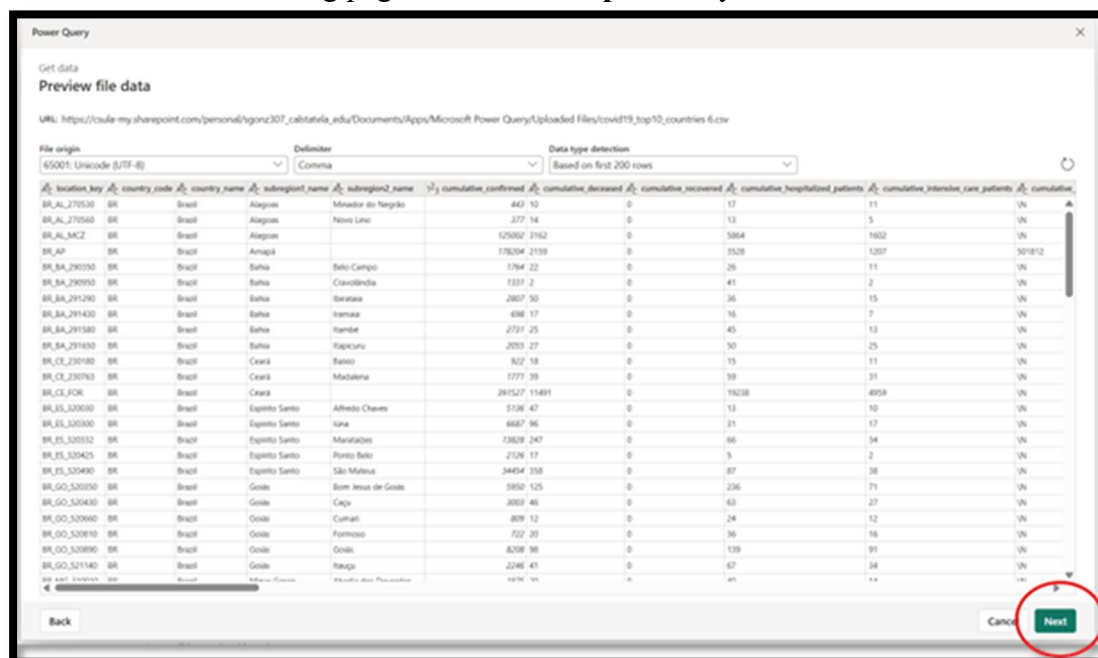Once you see the following page, select **csv**



Select the **Browse** button. Choose the file you want to use on the File Explorer, then hit **Next**
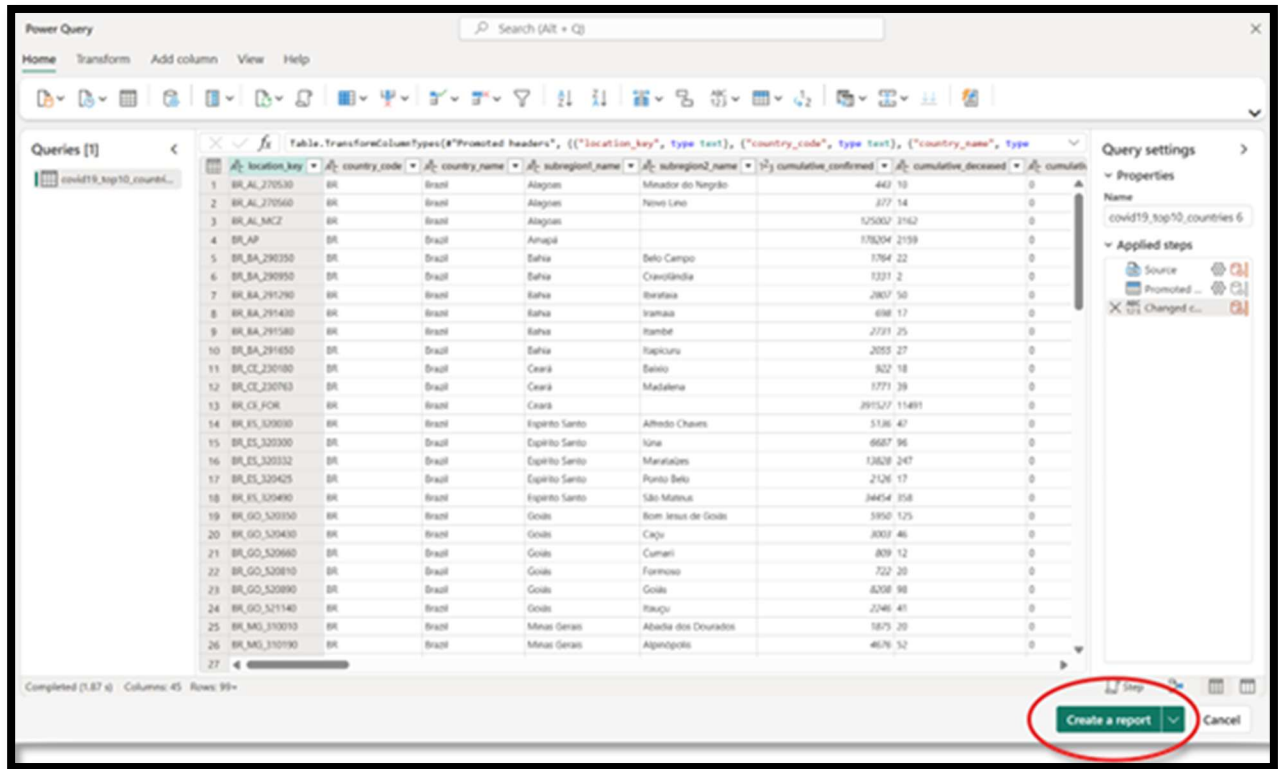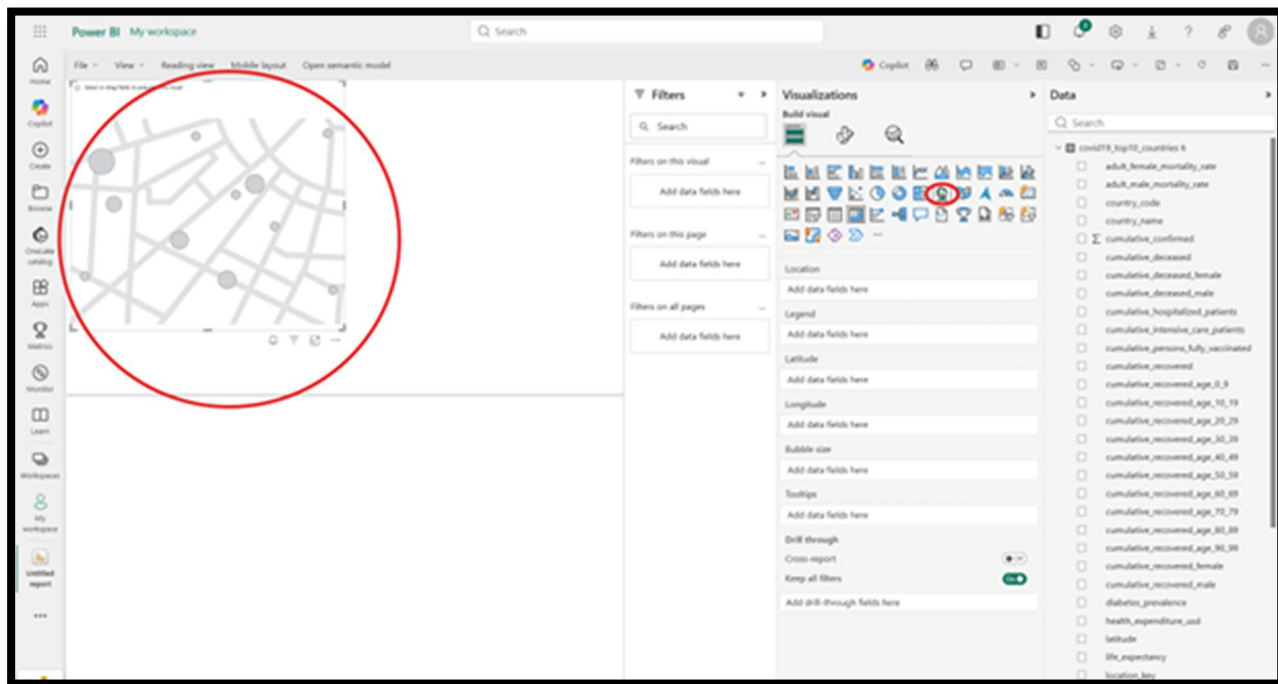
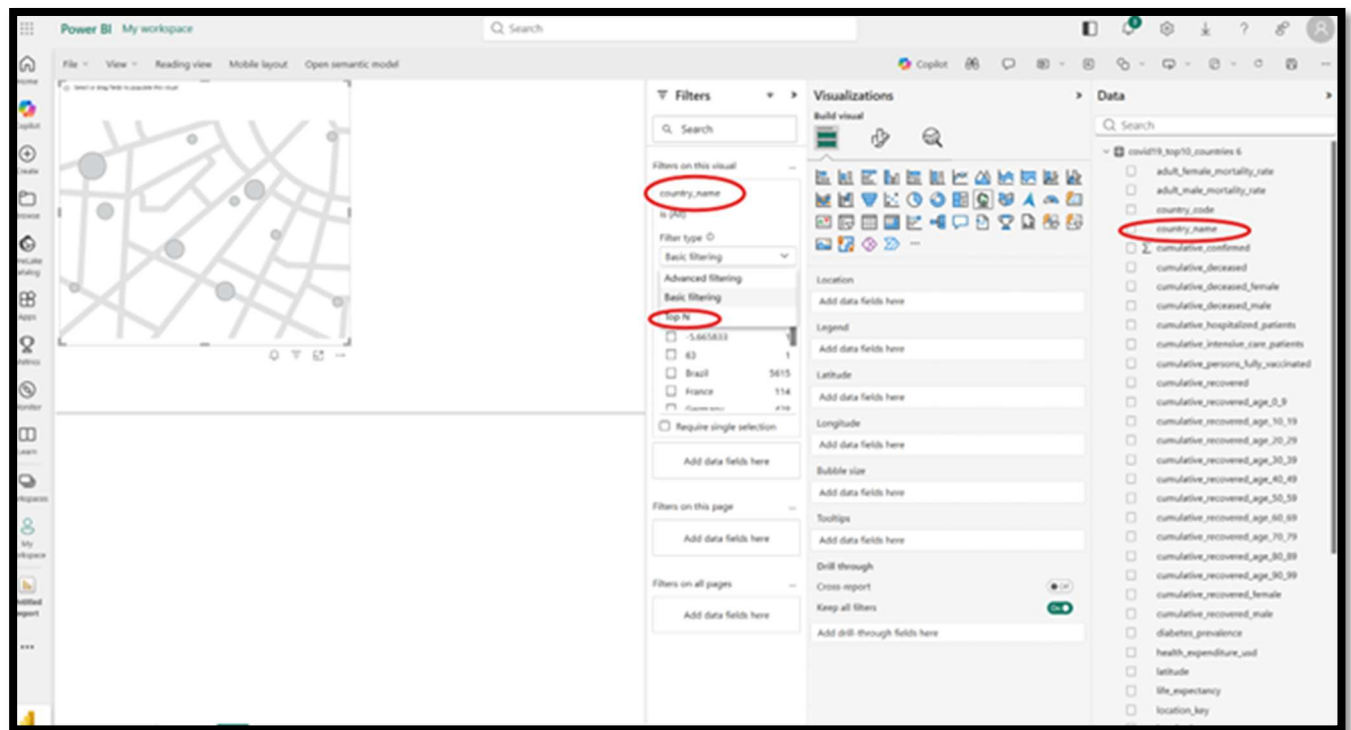You will see the following page in PowerBI to **preview** your dataset

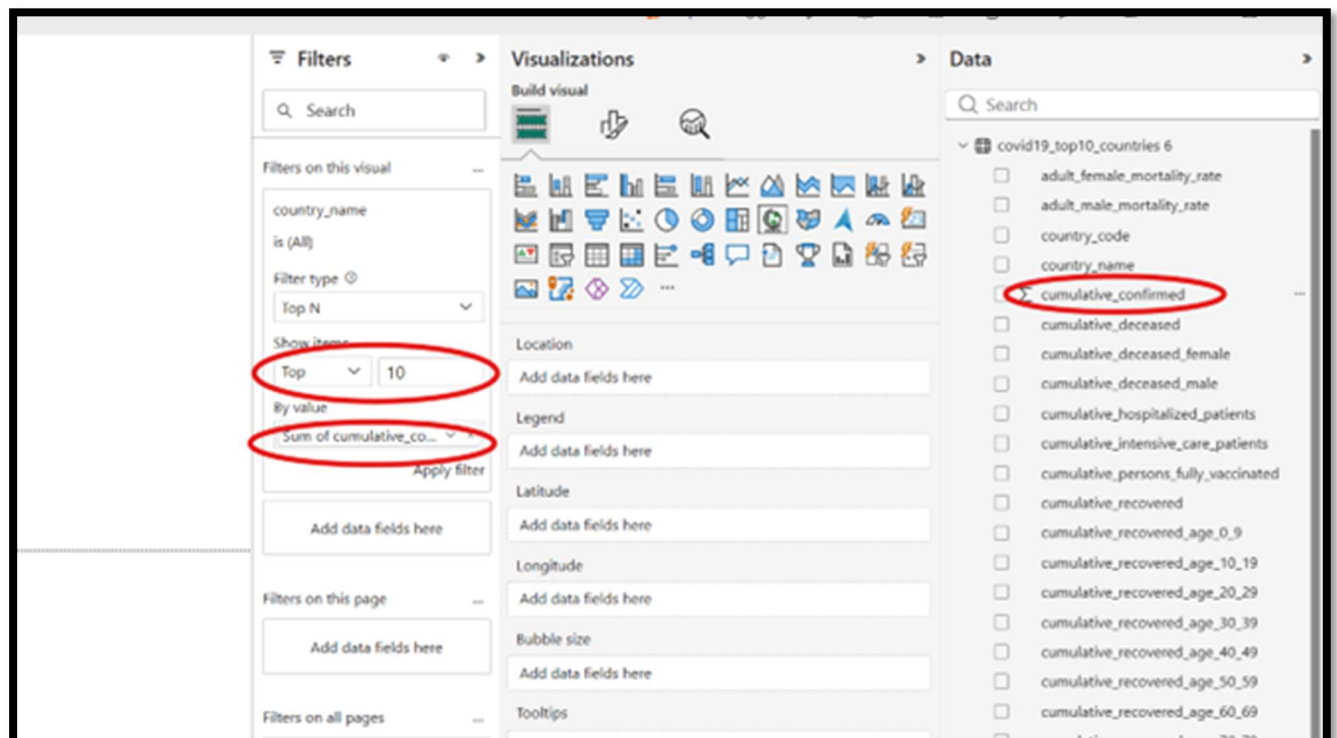On this page, select **Create a report -> Create**



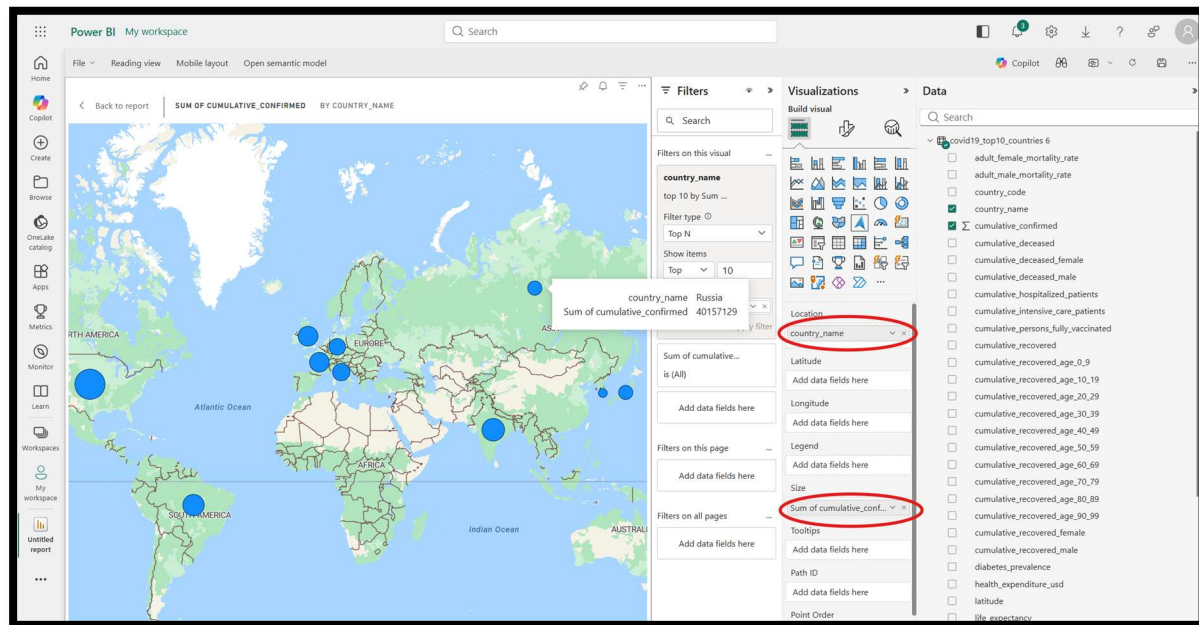Click on the **map** icon to begin setting up the visualization

Drag the **country_name** to the filters on this visual



On Filter type click Top N, then on show items click Top then next to it change it to **10** then drag the **cumulative_confirmed** to By value, then click **Apply filer**

Drag the **country_name** field to the **Location** area and the **cumulative_confirmed** field to the **Size** area



References

1. **COVID19-Data-Analysis: Final Project for CIS 4560. COVID-19 data analysis using HDFS Hive** - https://github.com/chica-94/COVID19-Data-Analysis