# Data Engineering Challenge

Proposed by Dell Technologies

Solved by Lucas Webber Molin

December 29, 2021 - January 2, 2022

Exercise 1 - Loading Data

First, I accessed the websites provided in the challenge documentation, and analyzed the JSON and CSV files to understand what I had in hand. Then I created a PostgreSQL database locally on my computer, and a table for each of the two files. I sent the CSV file to the countries_of_the_world_csv database table using the DBeaver application, where I also wrote the queries used in this challenge. For the JSON file, I created a Python project through the PyCharm IDE.

The script for creating the tables is attached under the name create_tables.sql in the SQL_Script folder. The data load is in the covid19_insert.sql and countries_of_the_world_csv_insert.sql files also in the SQL_Script folder.

Exercise 2 - Create a Pipeline

In the main.py script, attached to the Load_Script folder, I wrote a code to search the JSON file from the website and insert it into the covid19 database table. On the website, it says that the file is updated weekly, so to add only new data I implemented a condition to only insert records of dates longer than the longest date already stored in the database. To avoid duplicate records, I created a unique constraint. To automate the extraction, a task could be created in the Windows task scheduler to run weekly.

I noticed that between the two data sources there were the same countries with the name spelled differently. I used the query that is in the full_outer_join_countries.sql file attached to the SQL_Script folder to identify which countries they were. I performed the correction manually, looking up the countries on Google, using the covid data source name as the main one.

Exercise 3 - Create a View

The script for creating the view is in the view_ex3_create.sql file attached to the SQL_Script folder. And in the same folder is the view data sample in the view_ex3_insert.sql file.

Exercise 4 - Queries

The queries written in SQL for this exercise are in the queries_ex4.sql file attached to the SQL_Script folder, along with the performance analysis.

Exercise 5 - Data from other sources

As an additional data source, I chose tourism data for each country to identify correlations between the number of visitors and the number of COVID-19 cases. Seeing that small countries like Aruba, for example, had many cases for a few inhabitants, I wondered if there was a correlation.

Source: https://onestep4ward.com/most-visited-countries-in-the-world/
The data were extracted manually.

The data load is in the tourism_countries_insert.sql file, along with the tourism_countries_create.sql table creation file, in the SQL_Script folder.

Pearson's correlation calculation was done in Python, according to the correlation.py file attached to the Data_Analysis folder, using the query from the view_cases_visitors.sql file attached to the SQL_Script folder. The coefficient result was 0.60, indicating a moderate correlation.
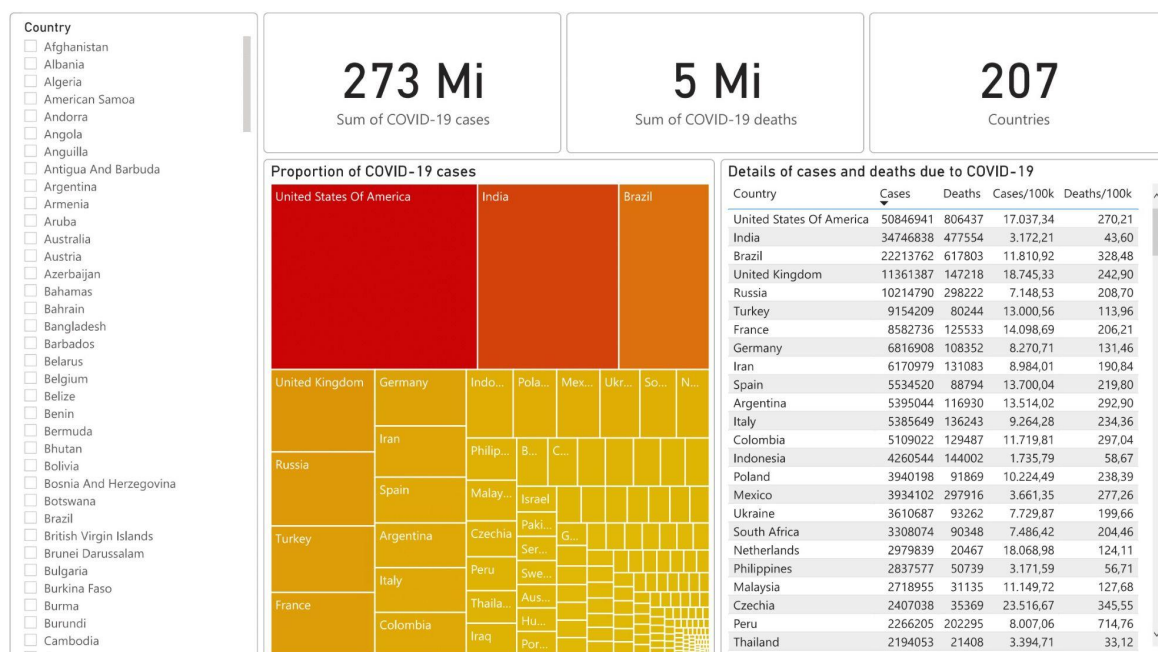
Exercise 6 - Report

To present the data, I chose to build a dashboard. The tool I use and master in my current job is protected by a confidentiality agreement, so I chose Power BI for this

challenge. But I'm a beginner at Power BI. The project file is attached in the Data_Analysis folder of this zip under the name dell2021.pbix

I built a simple dashboard with general data for total cases and deaths and for every 100,000 inhabitants. The dashboard contains a filter that allows you to select multiple countries; totalizers that show the sum of the number of cases, deaths and the number of filtered countries; a treemap showing the proportion of cases among the filtered countries; and a table with details for each country, which can be sorted by the chosen column (country, number of cases, number of deaths, case rate and death rate for every 100,000 inhabitants.

See the dashboard in the image below or in the attached dell2021.pbix file in the Data_Analysis folder.



A hypothesis that I wanted to test is the correlation between the number of cell phones and the number of cases, because it is common to distribute fake news about COVID-19 through social networks. The result of Pearson's correlation calculation was 0.55, which indicates a weak correlation. I imagine that more specific data would be needed to identify whether there is a correlation between the amount

of fake news spread and cases of COVID-19. The calculation was done in Python, according to the correlation.py file attached to the Data_Analysis folder, using the query from the view_cases_phones.sql file attached to the SQL_Script folder.

As requested, the conclusion of the data analysis is in the conclusion.txt text file in the Text_File_Conclusions folder.