

P4DS (ODL1) Assignment 3

Data Analysis Project

Notebook template design: Brandon Bennett (2020.11.03, revised 2021.03.02)

Analysing the Impact of Covid-19 on European Countries

Student: Loukas Tsourpilis
Email: od21lt@leeds.ac.uk

Project Plan

The Data (10 marks)

The data used for this project comes from ECDC, the European Agency for Disease Prevention and Control, an official agent of the European Union. There are two datasets that are used as part of this project, the official report of the ECDC daily covid cases/deaths across the EU/EUFA area and the vaccination data from the same region. These datasets can be found in the following links: [Link 1](#) [Link 2](#)

This dataset for the cases/deaths consists of 13 columns in total excluding the indexes. They are date, day, month, year, country, new cases for that date, new deaths for that date, country code, country geographical code, country population and continent. It is important to note that this dataset contains the data only from the 1st of March 2021 until the 28th of October 2021. Data before March were reported on a weekly basis, and are part of a different dataset. By observing the dataset, it becomes apparent that some piece of data do not offer any new information. For example we can identify the day/month/year by looking at the respective columns and as a result the date string is redundant. The same also applies to countries which have 3 different identifiers, while the continent column offers nothing as obviously this dataset concerns european countries only. Later when analysing this data we can trim some of these columns so the dataset becomes smaller and easier to handle with a computer programme.

The identifier for the vaccine data consists of 12 columns in total excluding the indexes. These are the year-week ISO (week of the year identifier), number of first doses administered, number of first doses refused, number of second doses administered, number of unknown doses administered, number of doses supplied to that country, region, country, population, target group, vaccine and denominator. The official and more detailed data dictionary can be found here. At first sight it is clear that the various european countries have a different reporting system. For instance, some countries report by region while some countries don't and just use the country identifier in the region column. It is important to keep that in mind when performing the analysis further down the road.

These datasets are constructed from ECDC using the official reports for covid cases/deaths of member states, as as a result they are the most accurate data available for the EU area. However there are a couple caveats that need to be noted. Firstly, the number of covid cases per country heavily depends on how many covid tests that specific country performs per day. Some countries may not invest heavily in detecting covid so as a result they have lower numbers in new cases. Secondly, since all the data mentioned above describes countries that differ in population significantly, we have to keep in mind that raw numbers offer minimal insight on the covid impact, with percentages and ratios being better descriptors of the situation.

Project Aim and Objectives (5 marks)

The Covid-19 pandemic has had a great impact on everyday life for the past two years. The general aim of this project is to take some very simple daily reported data for each European country, like daily covid cases, deaths and vaccine doses administered, and convert them to more meaningful information that can give a better insight on the situation so far across Europe. In more detail, the following three steps summarise the aim of this project:

- The fatality rate of the virus is measured for each European country. The fatality rate equals the total number of covid deaths divided by the number of covid cases and is used to be described in common terms how deadly the virus is. The fatality rate is a very useful indicator for two reasons. Firstly, the average fatality rate is a good indicator to determine the actual death rate of covid, and secondly it is important to see if the death rate varies significantly between different countries and why.
- The total cases per 100,000 population are calculated. The official report for each country usually contains the new cases for each day as a simple number. However in order to compare the overall situation between countries and see where covid is more prevalent, it is important to convert those numbers to a ratio per 100,000 population, which allows to compare numbers on equal terms. This ratio is calculated using the following formula (Events occurrence)/100,000/Population. This ratio is afterwards used to classify countries according to covid risk level. Although there is no standard framework to assess risk levels and each authority may use a different model, this project uses a framework developed by the Harvard Global Health Institute and can be found online [here](#).
- The correlation between vaccine doses administered and deaths from covid is found. A negative correlation is a good indication that the vaccines do indeed offer strong protection against death from covid. To calculate the correlation the following mathematical formula is used:

$$\rho_{X,Y} = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y}$$

For formula is used: $\text{cov}(X,Y)$, which is the standard method when using the pandas' built-in correlation function.

Specific Objective(s)

- Objective 1:** Measure the fatality rate of each country.
- Objective 2:** Calculate the case per population ratio and death rates according to covid risk level.
- Objective 3:** Find the correlation between vaccination and casualty rates.

System Design (5 marks)

Architecture

The following flowchart describes the flow of the programme. Each time a module is run we get a new stage which is described in the task section. Meanwhile the boxes represent the dataframes that exist in the memory. Some abbreviations are used. FR stands for fatality rate while DF for dataframe.

