# Exercise 02 - Statistic and Data Visualization

**Teacher: Le Ngoc Thanh - Tran Trung Kien - Nguyen Tien Huy**
**Student: Phan Vi Giai - 19127639**

## Task Table

| STT | Task | Separated task | % Completed |
|---|---|---|---|
| 1 | Crawling data | | 100% |
| 2 | Preprocessing data | | 100% |
| 3 | Visualizing data | Statistic and Conclusion | 100% |
| | | Visualize data with charts | 100% |

## Table of Contents

## Exercise 02 - Statistic and Data Visualization

# 1. Introduction

From about the end of 2019 until now, a plague has spread terribly around the world. Every day thousands of people are infected and tens to hundreds of people die. Organization Worldometer (www.worldometers.info) has collected statistical data from many sources and from many countries. Join daily reports to aggregate into a table in Figure .

In the website, the organization Worldmeter also makes graphs to show the visual change of the epidemic situation. However, we do not use it for the time being. We are mandated by the country to understand this data as between data fields that are related, what is the relationship, is there an anomaly in the data as the country report differs from the data synthesis, anomalies in data jumping, .... Apply knowledge of statistics and data visualization to make sense of existing data.



# 2. Data preparation

## 2.1 Crawling Data

### a. Approach

First, we need to get data for our statistics and visualization. But how do we get it ?
First of all, we access the *www.worldometers.info* and view the base HTML.
The strategy is we use Python and crawl data from this website.
Because the data we need is included in the general table tag in HTML structure.
It is so easy to get each data in the table cells.

We must crawl data at midnight for having almost full data of that day.

### b. *Algorithm*

The idea is looping through every <tr> tag, with each <tr> tag we crawl each <td> tag to the end of table cells.





```python
# Main crawling
def crawl(self):
    print('--------------------- Running Crawling Job')
    sleep(1)
    tabs ={
        'https://www.worldometers.info/coronavirus/#main_table': ['nav-today']
    }
    for tab in tabs.keys():
        self.browser.get(tab)
        sleep(1)
        selectors = tabs[tab]
        for selector in selectors:
            a = self.browser.find_element(By.XPATH, f"//a[contains(@href, '#{selector}')]")
            self.browser.execute_script("arguments[0].click();", a)
            sleep(1)
            self.get_info(selector)
    print('--------------------- Done Crawling Job')
    self.browser.close()
```

Our main function has a tab object which contains a key for link website and value is a table selector. The function which is crawling data is **get_info()**

```python
# Get information of specified selector
def get_info(self, selector):
    tab_panel = self.browser.find_element(By.ID, selector)
    table = tab_panel.find_element(By.TAG_NAME, 'table')
    # Get header attribute
    header = self.get_header('thead', table)
    # Body content
    bodys = table.find_elements(By.TAG_NAME, 'tbody')
    for body in bodys:
        # Specific world information
        if 'today' in selector:
            date = datetime.today().strftime('%d-%m-%Y')
        elif 'yesterday2' in selector:
            date = datetime.today() - timedelta(days=2)
            date = date.strftime('%d-%m-%Y')
        else:
            date = datetime.today() - timedelta(days=1)
            date = date.strftime('%d-%m-%Y')
        self.get_countries_info(body, header, f"data/{date}.csv")
    self.header = False
```

The idea of this function is:
- First, find the tab panel of argument selector which is passed
- Find the main table which contains our table cells
- Then, get the header from table element above

```html
▶<div id="main_table_controls_loading" style="width: 100px; display: none;">…
</div>
▶<div>…</div>
▶<div id="main_table_controls" style="position: relative;">…</div>
▼<div class="tab-content" id="nav-tabContent">
  ▼<div class="tab-pane active" id="nav-today" role="tabpanel" aria-labelledby=
  "nav-today-tab">
    ▼<div class="main_table_countries_div">
      ▶<div id="main_table_countries_today_wrapper" class="dataTables_wrapper f
      orm-inline dt-bootstrap no-footer">…</div>
      ▼<table id="main_table_countries_today" class="table table-bordered table
      -hover main_table_countries dataTable no-footer" style="width: 100%; margi
      n-top: 0px !important;"> == $0
        ▶<thead>…</thead>
        ▶<tbody>…</tbody>
        ▶<tbody class="body_continents">…</tbody>
        ▶<tbody class="total_row_body body_world">…</tbody>
      </table>
    </div>
  </div>
  ▶<div class="tab-pane " id="nav-yesterday" role="tabpanel" aria-labelledby="n
  av-yesterday-tab">…</div>
  ▶<div class="tab-pane " id="nav-yesterday2" role="tabpanel" aria-labelledby=
  "nav-yesterday2-tab">…</div>
</div>
▶<div style="font-size:13px; margin-top:10px; padding-bottom:10px">…</div>
▶<div style="clear:both; font-size:13px; margin-top:25px; padding-bottom:45px">
```

- Because the table contain some <tbody> tags, so that we must find all <tbody> tags and loop through for collecting each table cells
- Then calling the **get_countries_info()** for collection each data cell in the table

```python
# Get all countries cases information
def get_countries_info(self, selector, header, file_name):
    trs = selector.find_elements(By.TAG_NAME, 'tr')
    res = []
    for tr in trs:
        res.append(self.get_country_info(tr))
        if len(res) > 40:
            res = self.remove_empty_rows(res)
            self.write_csv(file_name, res, header)
            res.clear()
    if len(res) > 0:
        res = self.remove_empty_rows(res)
        self.write_csv(file_name, res, header)
```

The idea of this function is:
- First, find all <tr> tags and store in trs variable



- Loop each <tr> tag and append data to res variable
- Final, we write result to target file with strategy is saving batch size

### c. *How to run*

First, we open the folder which contain our source code

File **virus_crawl.py** is the main file for crawling data from the website.
Then open the terminal in this folder and run the command **python virus_crawl.py**



Finally, we check the result in the **data folder** which is raw data from the website we parsed.

## 2.2 Preprocessing Data

### a. Approach

After crawling data from the above step, we then check our raw data to find whether this has a fault type, missing data, or something else… Just open our csv file.



Glance at our raw data. We just find some cells are missing and some cells have N/A value.

| 23 | 21 | Malaysia | 2632782 | | 30425 | | 2537204 | | 65153 | 509 | 79899 | 923 | 37735175 | 1145179 | 32951341 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 24 | 22 | Peru | 2236351 | | 201176 | N/A | N/A | N/A | | 971 | 66518 | 5984 | 20314089 | 604225 | 33620094 |

Look at the row 23 and 24 which have problems.

### b. Algorithm

The idea for preprocessing data is to replace all missing values and N/A value with empty

value.

Our algorithm is to loop through each row and replace missing values or N/A values with empty values for easily counting missing values in the later step. Then save the preprocessed data to a new folder.

```python
# Main function for preprocessing raw data
def preprocess(self):
    print('-------------------- Running Preprocessing Job')
    for i in range(len(self.data)):
        for j in range(len(self.data[i])):
            if self.data[i][j] == '' or self.data[i][j] == 'N/A':
                self.data[i][j] = ''
            else:
                if self.data[i][j][0] == '+':
                    self.data[i][j] = self.data[i][j][1:]
                if self.data[i][j][0].isdigit() and i != 1 and i != len(self.data) - 1:
                    self.data[i][j] = np.int64(''.join(self.data[i][j].split(',')))
    self.header = self.data[0]
    self.data = self.data[2: len(self.data) - 1]
    self.data = [row[:-1] for row in self.data]
    self.write_csv(f'preprocess-data/preprocessed-{self.file_name}')
    print('-------------------- Done Preprocessing Job')
```

To run this preprocess task, the constraint is that we must has raw data first by running a crawling data task.

The **preprocess function** is the main function, which do some specific task:

- Running out each row and replace each cell has empty or N/A value with empty value
- Parse header field with the first row of raw data
- Finally, save our preprocessed data with 'preprocessed' prefix of file name

```python
# Main function for preprocessing raw data
def preprocess(self):
    print('-------------------- Running Preprocessing Job')
    for i in range(len(self.data)):
        for j in range(len(self.data[i])):
            if self.data[i][j] == '' or self.data[i][j] == 'N/A':
                self.data[i][j] = ''
            else:
                if self.data[i][j][0] == '+':
                    self.data[i][j] = self.data[i][j][1:]
                if self.data[i][j][0].isdigit() and i != 1 and i != len(self.data) - 1:
                    self.data[i][j] = np.int64(''.join(self.data[i][j].split(',')))
    self.header = self.data[0]
    self.data = self.data[2: len(self.data) - 1]
    self.data = [row[:-1] for row in self.data]
    self.write_csv(f'preprocess-data/preprocessed-{self.file_name}')
    print('-------------------- Done Preprocessing Job')
```

- In the field 'New Cases' and 'New Deaths' and 'New Recovered', data has '+' in the prefix indicating some cases added. So that we must remove a '+' in the prefix
- All the number has format ',' so we must remove the ',' and cast the type to int

**Advantage** of this approach is:

- Easy
- Fast
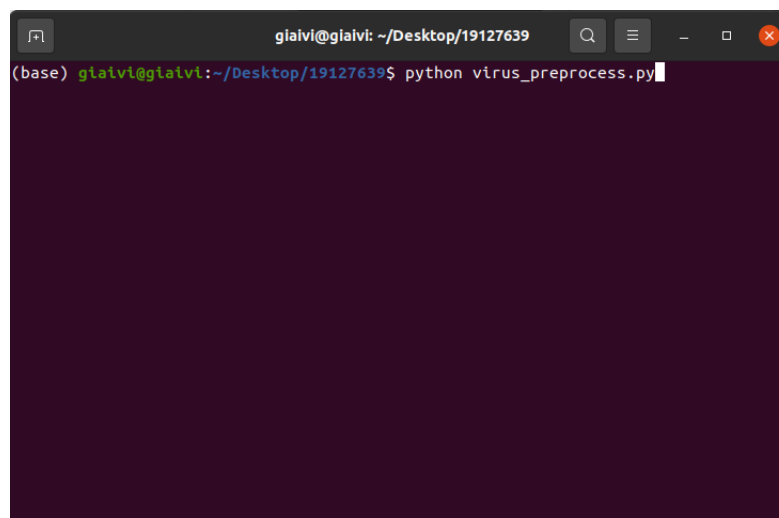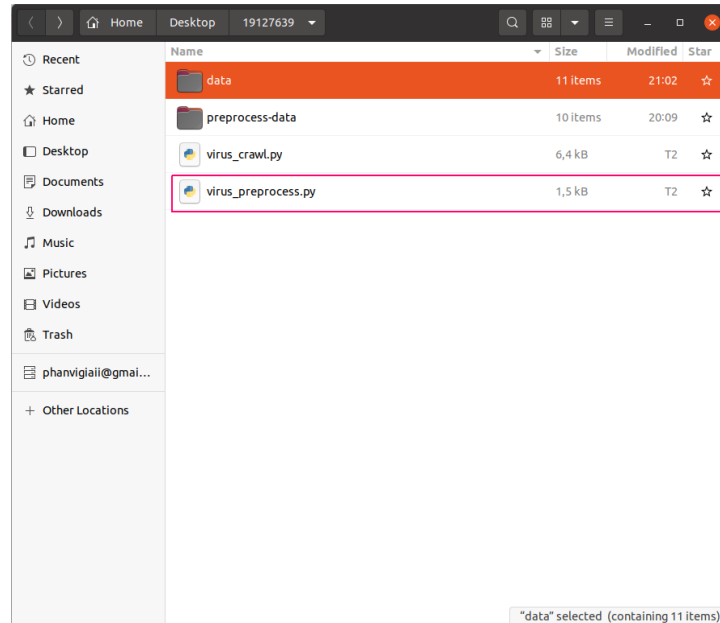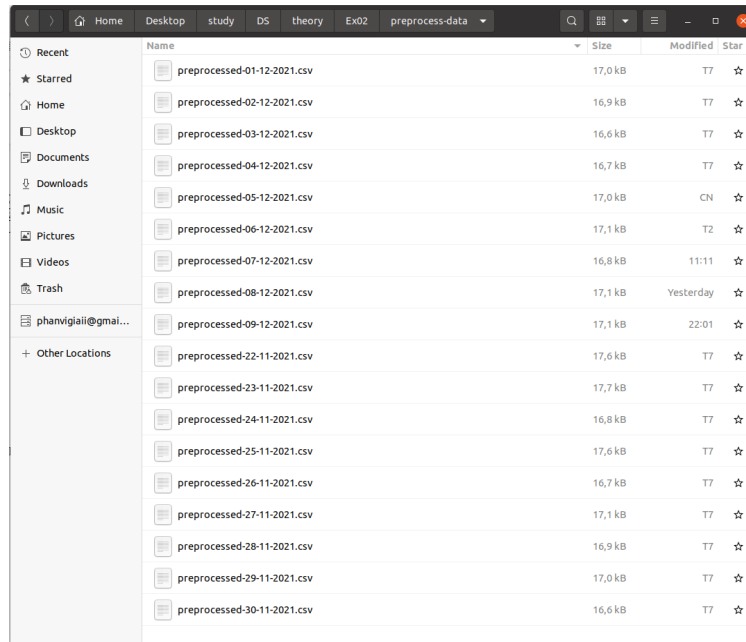- No performance for calculating

## c. *How to run*

To run this preprocess data job, we need to run crawl data first.
First, run crawl data to get our raw data.
Then open a folder which contains our source code.





Run command line **python virus_preprocess.py**

Then check the result in **folder preprocess-data**



missing data or N/A value is replaced with nan value

## 2.3    Automatic crawl and preprocess data with Apache AirFlow

### a.  *Approach*

In general, we manually run two command line for crawling and preprocessing data:

- python virus_crawl.py
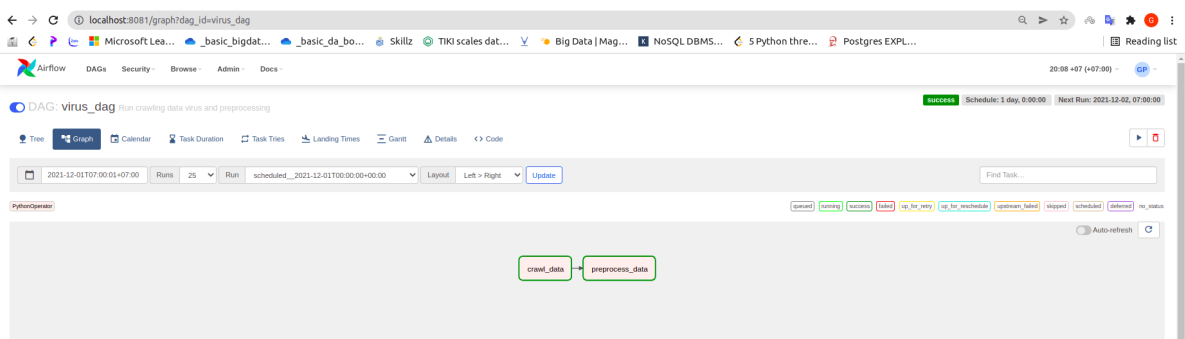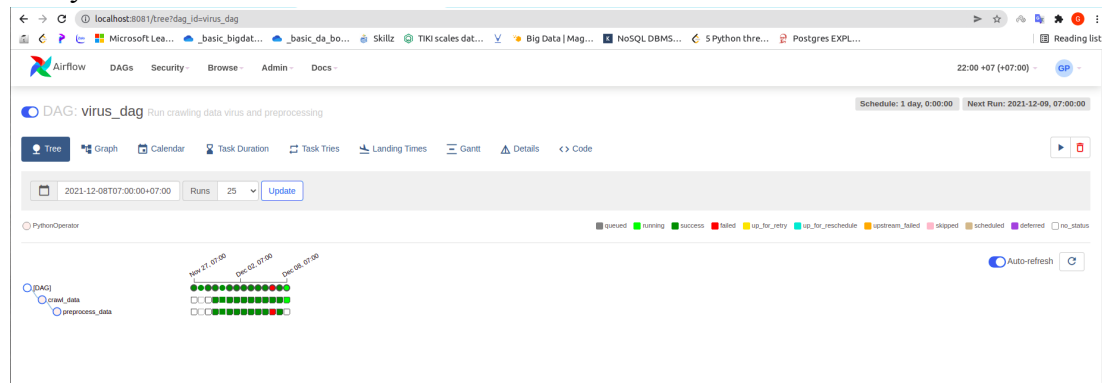- python virus_preprocessing.py

In this project, we must run two command line above ***daily*** for collecting data from the website: https://www.worldometers.info/coronavirus/

But, can we do those tasks in an automatic way ? Yes, we can do it.

What is the AirFlow ? Airflow is a platform to programmatically author, schedule and monitor workflows or data pipelines.

Using the AirFlow, we can automatically schedule our tasks. In this project, we build one data pipeline for crawling and then preprocess our data **daily at midnight** to collect almost full data of that day.
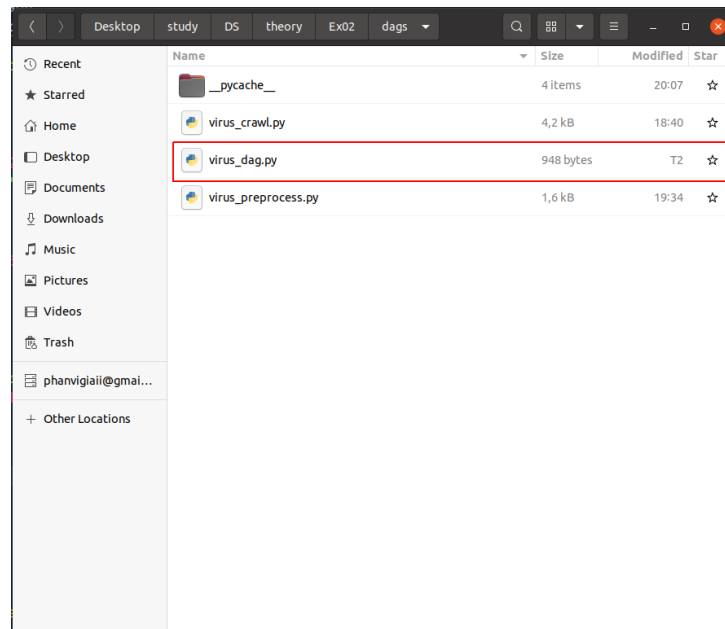




AirFlow dashboard

## b. Algorithm

To build the data pipeline, we need the **dag.**

What is **dag** ? **Dag** contains these attributes:

- Workflow as a Directed Acyclic Graph (DAG) with multiple tasks which can be executed independently.
- Airflow DAGs are composed of Tasks.

First, we create a **dag file** named **virus_dag.py**
Following these step in **Apache AirFlow** website
**Step 1:** Importing necessary modules

```python
from airflow import DAG
from datetime import timedelta, datetime
from airflow.utils.dates import days_ago
from airflow.operators.python_operator import PythonOperator
from virus_crawl import run_etl_crawl
from virus_preprocess import run_etl_preprocess
```

**Step 2:** Define default argument

```python
default_args = {
    'owner': 'GiaiVi',
    'depends_on_past': False,
    'start_date': datetime(2021, 11, 27),
    'email': ['phanvigiaii@gmail..com'],
    'email_on_failure': False,
    'email_on_retry': False,
    'retries': 1,
    'retry_delay': timedelta(minutes=5),
}
```

- Owner is me
- Start date: from 27-11-2021
- Retries: 1 time
- Retry delay: 5 minutes after failure

**Step 3:** Instantiate a DAG

```python
dag = DAG(
    'virus_dag',
    default_args=default_args,
    description='Run crawling data virus and preprocessing',
    schedule_interval=timedelta(days=1)
)
```

- dag id: virus_dag
- Describe the dag with description field
- schedule_interval: 1 day (mean daily)

**Step 4:** Specify our tasks

```python
run_crawl = PythonOperator(
    task_id = 'crawl_data',
    python_callable = run_etl_crawl,
    dag = dag
)

run_process = PythonOperator(
    task_id = 'preprocess_data',
    python_callable = run_etl_preprocess,
    dag = dag
)
```

- We have two task: crawling data and preprocessing data
- Use PythonOperator for calling python function
- task_id is our task name: crawl_data and preprocess_data
- python_callable is passed with our imported function

**Step 5:** Setting up dependencies

```python
run_crawl >> run_preprocess
```

- We specify **run_process task** depend on **run_crawl**
- Because we need raw data for preprocessing

## c. *How to run*

First, install AirFlow from official website

Activate AirFlow environment with command: **source airflow-venv/bin/activate**



Running airflow server with command: **airflow webserver -p 8081**
Run server with port 8081

Open another terminal, activate airflow environment like above step
Then, running airflow scheduler server with command: **airflow scheduler**

Finally, open web browser and access to **http://localhost:8081/**



AirFlow dashboard with **virus_dag** created above

Click on our dag and we have a control panel for manipulating tasks and so on...



Our scheduler

View dependable tasks

That is all, with AirFlow we can automatically run our job. But in this level, we don't have a server for deployment to running AirFlow server and scheduler 24/7.

Every day, the only thing we just do is run one command: **airflow scheduler**. The scheduler starts and automatically runs our job.



View a result in **folder data** which contains raw data and **folder preprocess-data** which contains preprocessed data.

# 3.    Statistic and Data Visualization

## 3.1    Expectations

After finishing to collect data from the website, what do we expect from our data ? In my expectations:
- We will get some columns depending on others like covariate, inverse
- We expect main column will be in standard deviation
- Percentage losing of data very small
- Type of columns are correct and reasonable

## 3.2    Statistic and Conclusion

### a. Statistic

First, we read the **"preprocessed-04-12-2021.csv"** file in the preprocess-data folder. We have the result below:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 224 entries, 0 to 223
Data columns (total 15 columns):
 #   Column             Non-Null Count   Dtype
---  ------             --------------   -----
 0   #                  224 non-null     int64
 1   Country, Other     224 non-null     object
 2   Total Cases        224 non-null     int64
 3   New Cases          18 non-null      float64
 4   Total Deaths       212 non-null     float64
 5   New Deaths         15 non-null      float64
 6   Total Recovered    216 non-null     float64
 7   New Recovered      17 non-null      float64
 8   Active Cases       216 non-null     float64
 9   Serious, Critical  157 non-null     float64
 10  Tot Cases/ 1M pop  222 non-null     float64
 11  Deaths/ 1M pop     210 non-null     float64
 12  Total Tests        209 non-null     float64
 13  Tests/ 1M pop      209 non-null     float64
 14  Population         222 non-null     float64
```

In total, data has 15 columns and 224 rows.

Each column has the meaning:

- "#" mean id of the row
- "Country, Other" mean the country which is reported about virus
- "Total Cases" mean total cases are calculated from the way it's always been.
- "New Cases" mean cases people are inflected in that day
- "Total Deaths" mean total deaths from that country until that day
- "New Deaths" mean new deaths happen in that day
- "Total Recovered" mean total cases recovered until that day
- "New Recovered" mean new cases recovered in that day
- "Active Cases" mean cases is inflected in that day which is calculated by "Total Cases" minus "Total Deaths"
- "Serious, Critical" mean some cases in danger, severe illness
- "Tot Cases/ 1M pop" mean total cases but is calculated by 1 million population
- "Deaths/ 1M pop" mean total deaths but is calculated by 1 million population
- "Total Tests" mean total tests which country tested virus for people until that day
- "Tests/ 1M pop" mean total tests but is calculated by 1 million population
- "Population" mean total number of people in the country

```
In [6]:  # Check if data has duplicated by checking the id
         df['#'].duplicated().any()

Out[6]: False
```

In data, there is no duplicated row by using built-in function of DataFrame to check it and return False

```
In [9]: df.dtypes

Out[9]: #                    int64
        Country, Other       object
        Total Cases          int64
        New Cases            float64
        Total Deaths         float64
        New Deaths           float64
        Total Recovered      float64
        New Recovered        float64
        Active Cases         float64
        Serious, Critical    float64
        Tot Cases/ 1M pop    float64
        Deaths/ 1M pop       float64
        Total Tests          float64
        Tests/ 1M pop        float64
        Population           float64
        dtype: object
```

Column "#", "Total Cases" has type int64 for numerical value
Column "Country, Other" has type object for string value
The rest have type float64 for numerical value

```
ColName             Miss(%) NumDifVals SomeVals
#                   0.000   224        {'10', '19', '35', '87', '54', '11...
Country, Other      0.000   224        {'Egypt', 'Kyrgyzstan', 'Nicaragua...
Total Cases         0.000   223        {'16000', '10', '5701', '91993', '...
New Cases           91.964  18         {'31', '138', '17659', '1702', '14...
Total Deaths        5.357   205        {'718', '35', '880', '3108', '387'...
New Deaths          93.304  12         {'10', '3', '70', '47', '20', '37'...
Total Recovered     3.571   213        {'860595', '84834', '34322', '2154...
New Recovered       92.411  17         {'22765', '34900', '160', '5666', ...
Active Cases        3.571   199        {'10', '72323', '782', '88', '828'...
Serious, Critical   29.911  99         {'10', '13714', '220', '110', '47'...
Tot Cases/ 1M pop   0.893   221        {'1901', '19', '13720', '80737', '...
Deaths/ 1M pop      6.250   198        {'10', '782', '828', '1766', '1181...
Total Tests         6.696   209        {'11925712', '37549349', '3580510'...
Tests/ 1M pop       6.696   209        {'6552', '1560323', '1056575', '12...
Population          0.893   222        {'11994233', '18023312', '30517', ...
```

The image above summary missing and different values which is calculated from our data
Base on the result:
- The columns "#", "Country, Other" and "Total Cases" have 0% missing mean that those attributes has no missing value and two divide three has unique values but not the column "Total Cases" has duplicated in two last row
- The columns "New Cases" and "New Deaths" are the largest missing percentages. The Ministry of Health can not record new cases and new deaths everyday because of the exploding infection rate in the country
- The remain has normal value

## b. *Conclusion*

Through some statistical information, we can recognize that the inflection of CoronaVirus has exploded all over the world.

The new types of CoronaVirus has been raise, quote from one paper:

"A new COVID-19 variant detected in South Africa with a high number of mutations has raised concerns among scientists and triggered travel restrictions by a number of countries amid fears of coronavirus transmissions. The National Institute for Communicable Diseases (NICD) said 22 positive cases of the new variant have been recorded in the country following genomic sequencing. News of the announcement broke on Thursday…"

That is the alarm for all people in the world of the dangerous dead if the government has no countermeasures.

***The survival of people depends on their consciousness and the sacrifices of doctors around the world.***
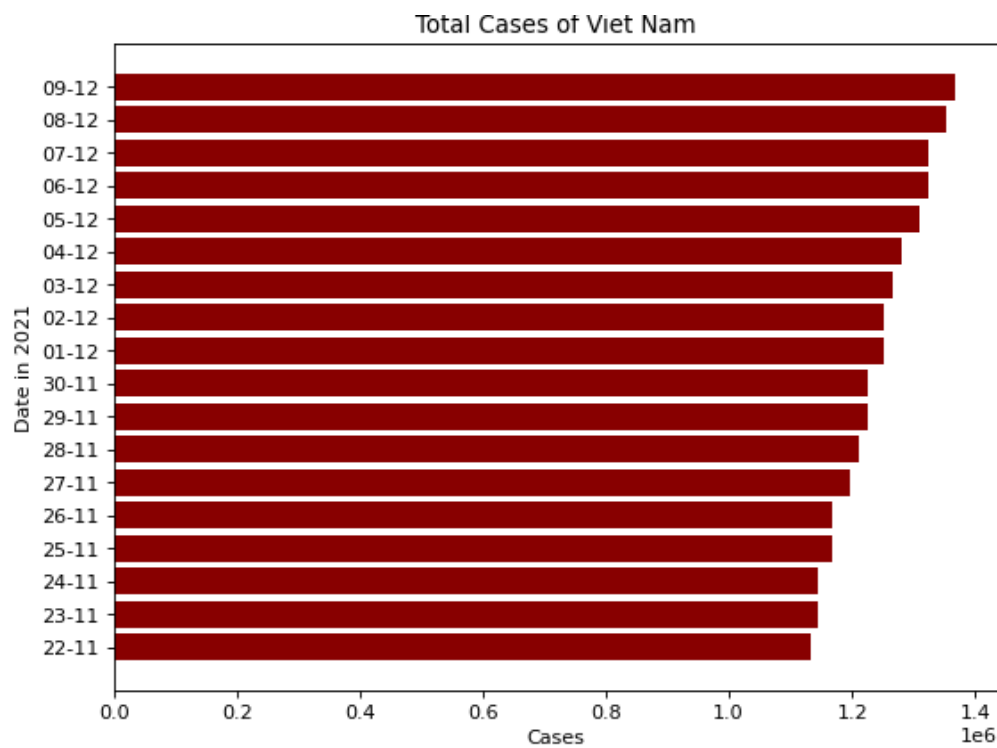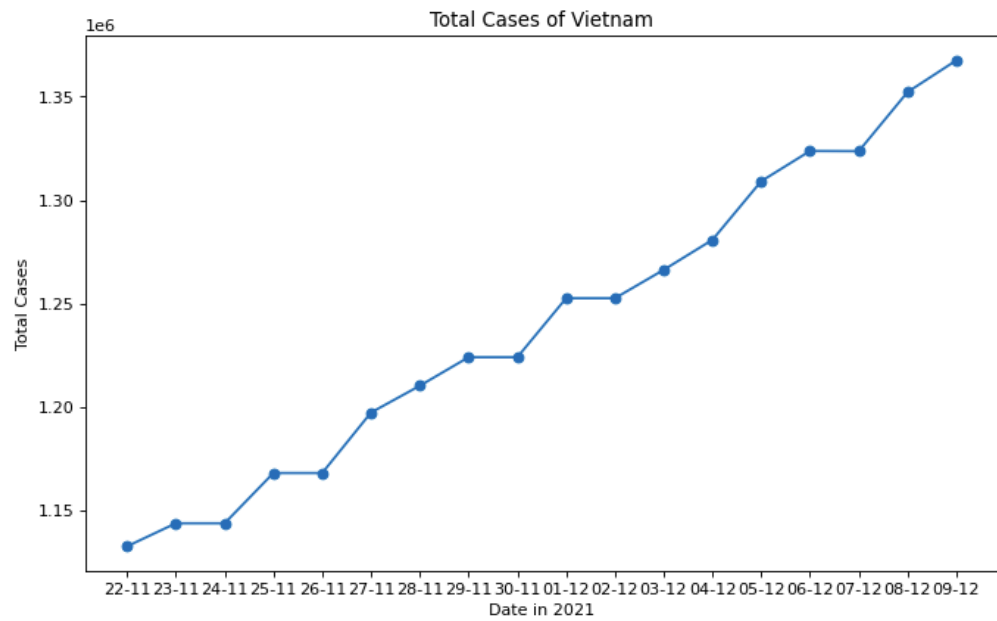
## 3.3 Visualization

### a. *Choose column for visualizing*

Because "#" is a column specifying Id for each row  and "Country, Other" is a column describing data of where. So that we choose the remaining column for the variable of the visualization chart.
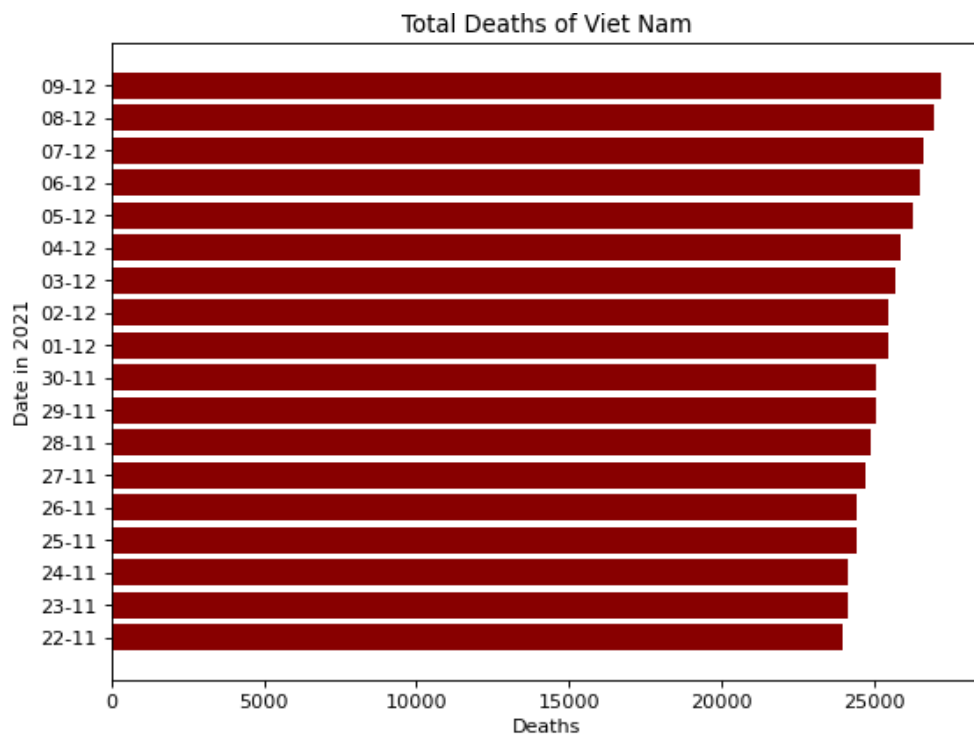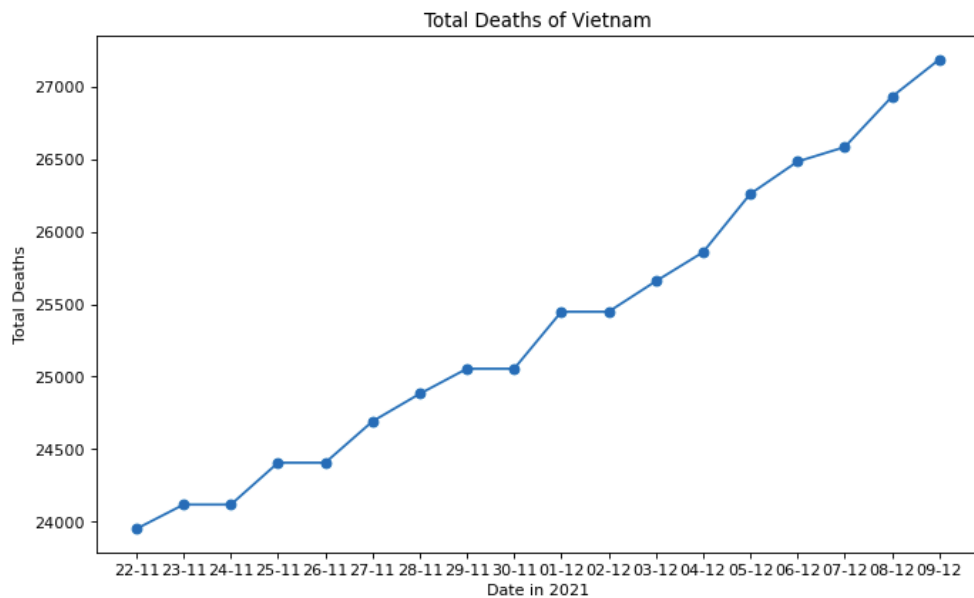
### b. *Comparison Plots*

First, we visualize each column of data with a relation date. To know how the trend of cases increases or decreases until now ? In this part, we use **line charts** and **bar charts** for visualizing data because three columns "Total Cases", "Total Recovered" and "Total Deaths" are numerical types and we want to see how the data change overtime.
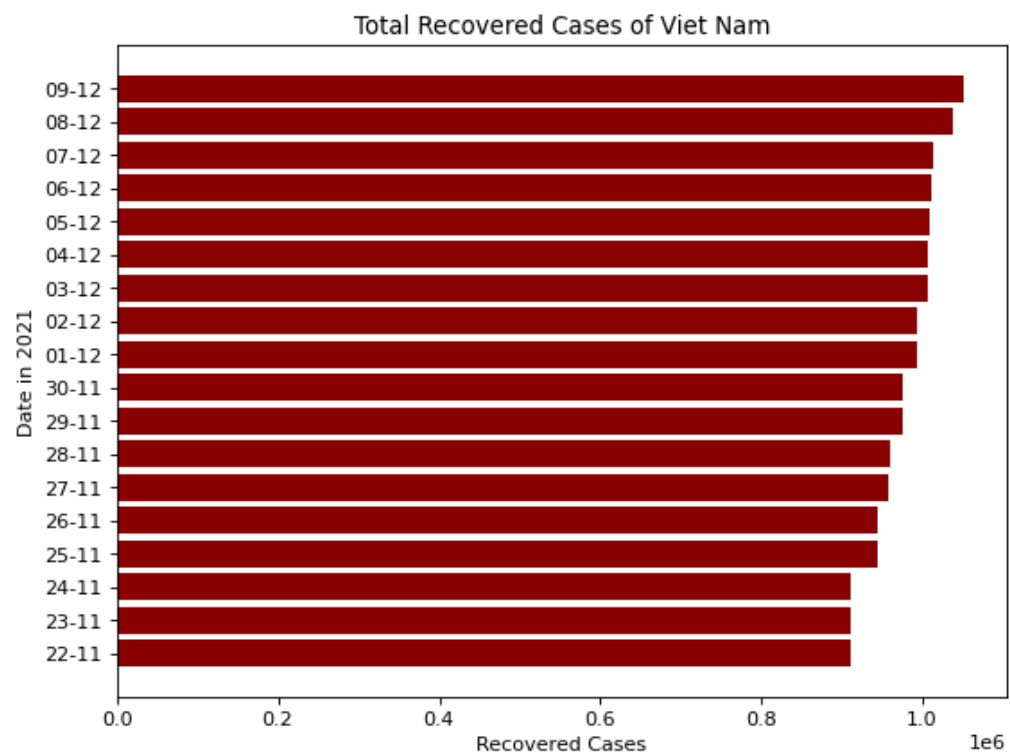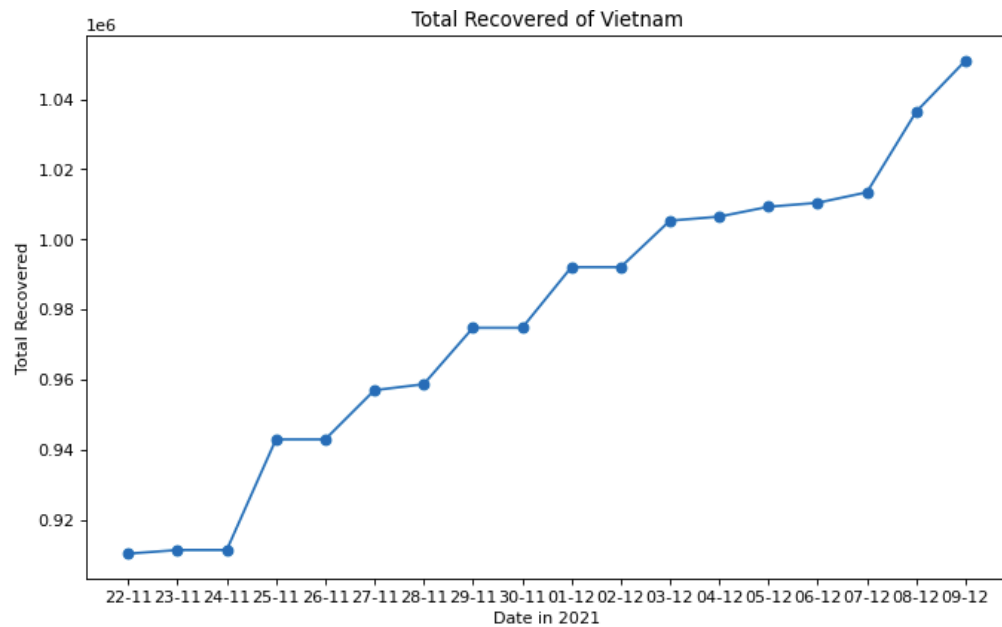
Total Cases of Vietnam



Total Cases of Viet Nam

Total Deaths of Vietnam



Total Deaths of Viet Nam

| Introduction To Data Science | Date: 30/11/2021 |
|---|---|
| Exercise 02 - Data visualization | Author: Phan Vi Giai - 19127639 |
| | |



Total Recovered of Vietnam



Total Recovered Cases of Viet Nam

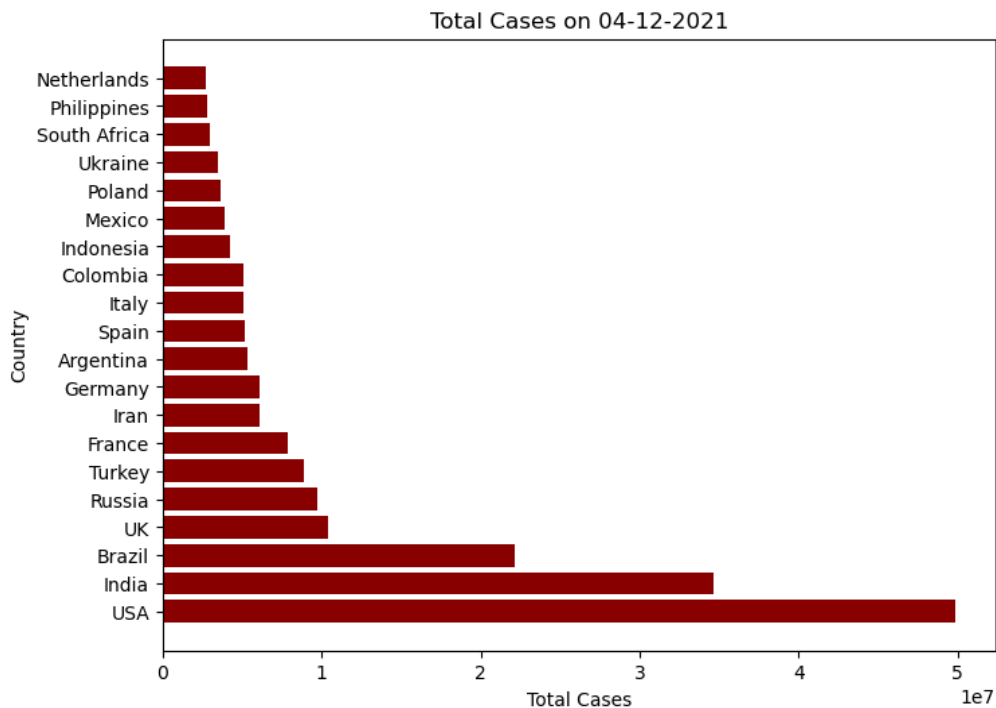Above charts indicate how the trend of CoronaVirus happening in VietNam through 22-11-2021 until 09-12-2021.

Total number of people infected by Virus is increasing cause total deaths also increase. But somehow, one positive thing is that the number of people recovering from the disease is increasing day by day.
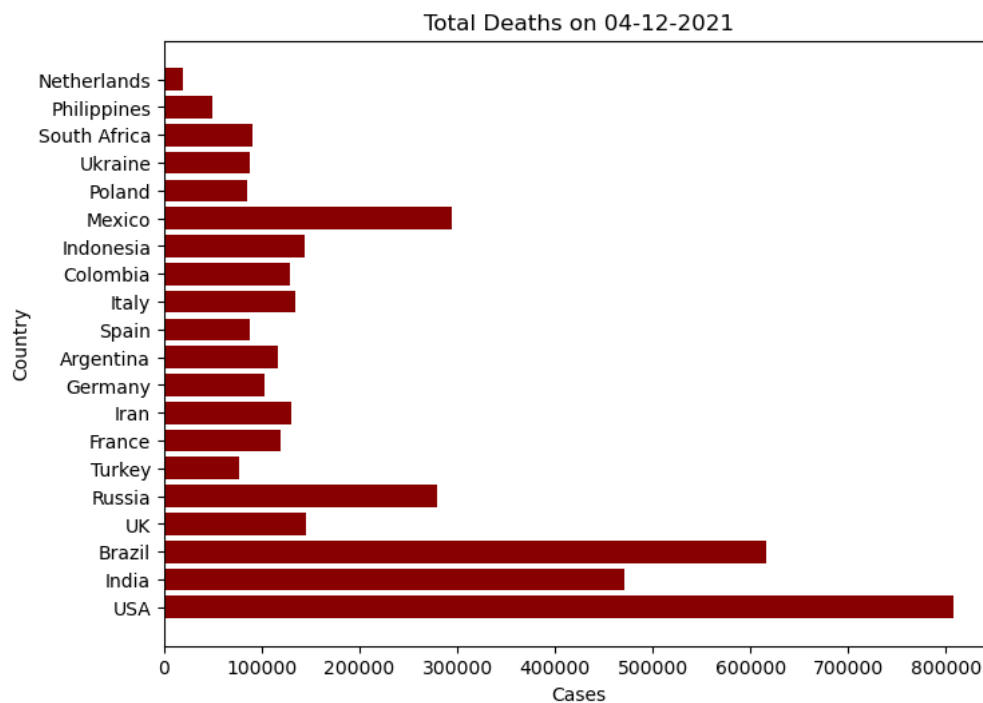


This chart shows a comparison between total inflected cases and total deaths in Viet Nam. "Total Deaths" occupies a small part of the whole. But let's not be subjective about this

dangerous disease. An example is the terrible virus disease in India and the US.



*Bar chart show top 20 country leading the number of Total Cases in the world*

*Bar chart show top 20 country leading the number of Total Deaths in the world*



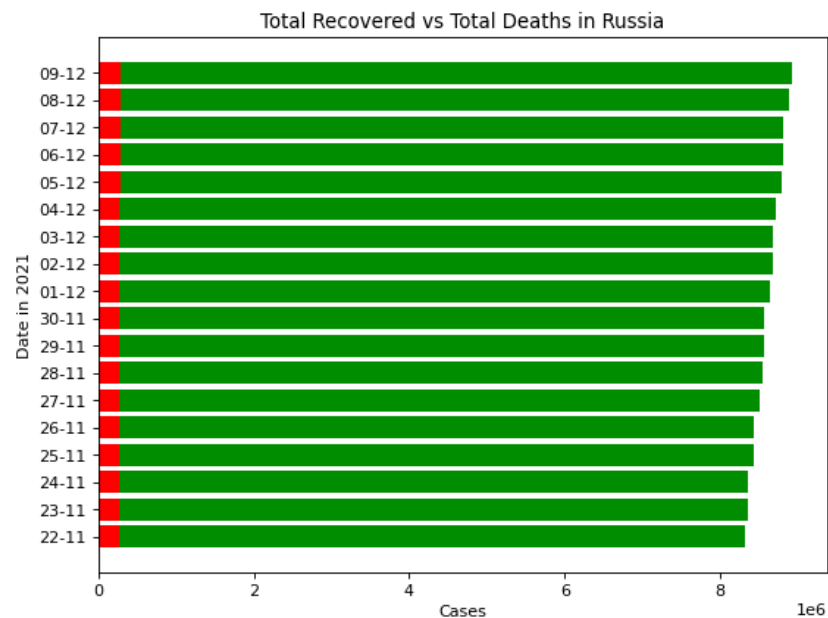*Bar chart show top 20 country leading the number of Total Recovered in the world*

The top of the three charts above is the US. The charts show that the disease covers the whole US because of the subjectivity and contempt for the disease since it has not yet spreaded. The ranked 2nd and 3rd are poor and backward countries are Brazil and India. The disease also spread and took many lives.

### c. *Composition Plots*

In this part, we use **Stacked Bar Chart** to show how each variable is divided into sub-categories and the proportion of the sub-category. We use "Total Recovered" and "Total Deaths" columns for visualizing because the two column meeting condition is a numeric attribute.
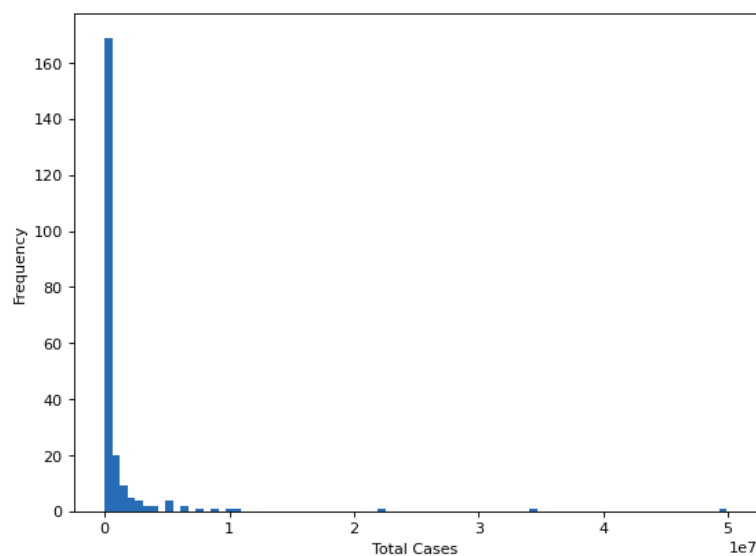
Data is collected in **Russia**.

"Total Deaths" variable occupies less part than the "Total Recovered" variable.
Russia is also a typical example like Vietnam.
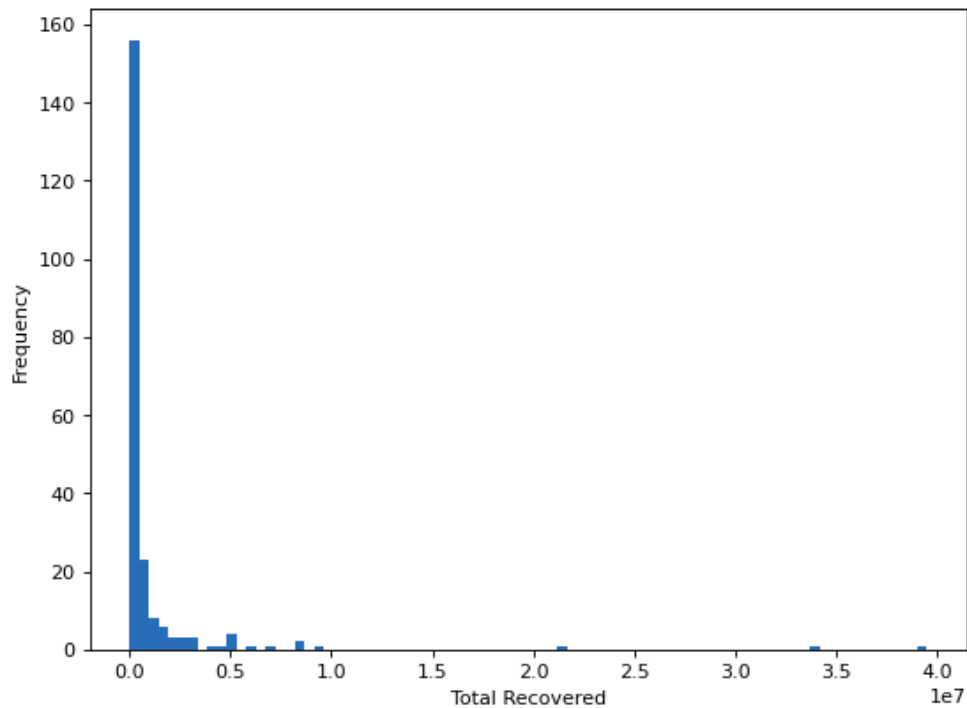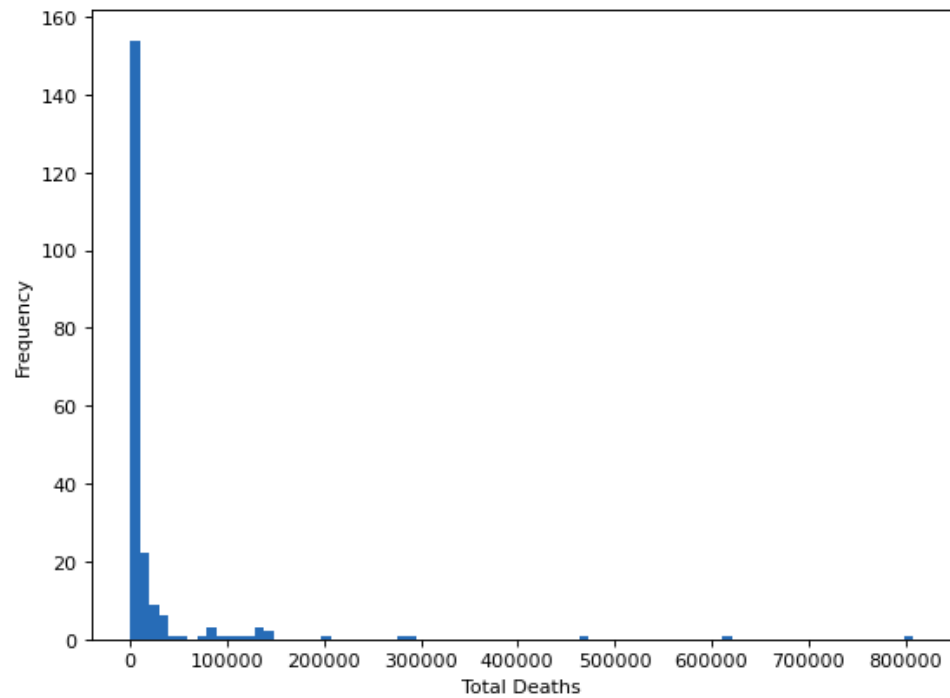
## d. *Distribution Plots*

We use two types of Distribution Plots: **Histogram** and **Boxplot**.

We want to know more about the distribution of "Total Cases", "Total Deaths", "Total Recovered" and its frequency. Also, see some statistical measurements like: Median, Mean, Upper Quartile, Lower Quartile,...

Data we use in this example is on **04-12-2021** for **Histogram** and data we use for **Boxplot** is of **Russia.**
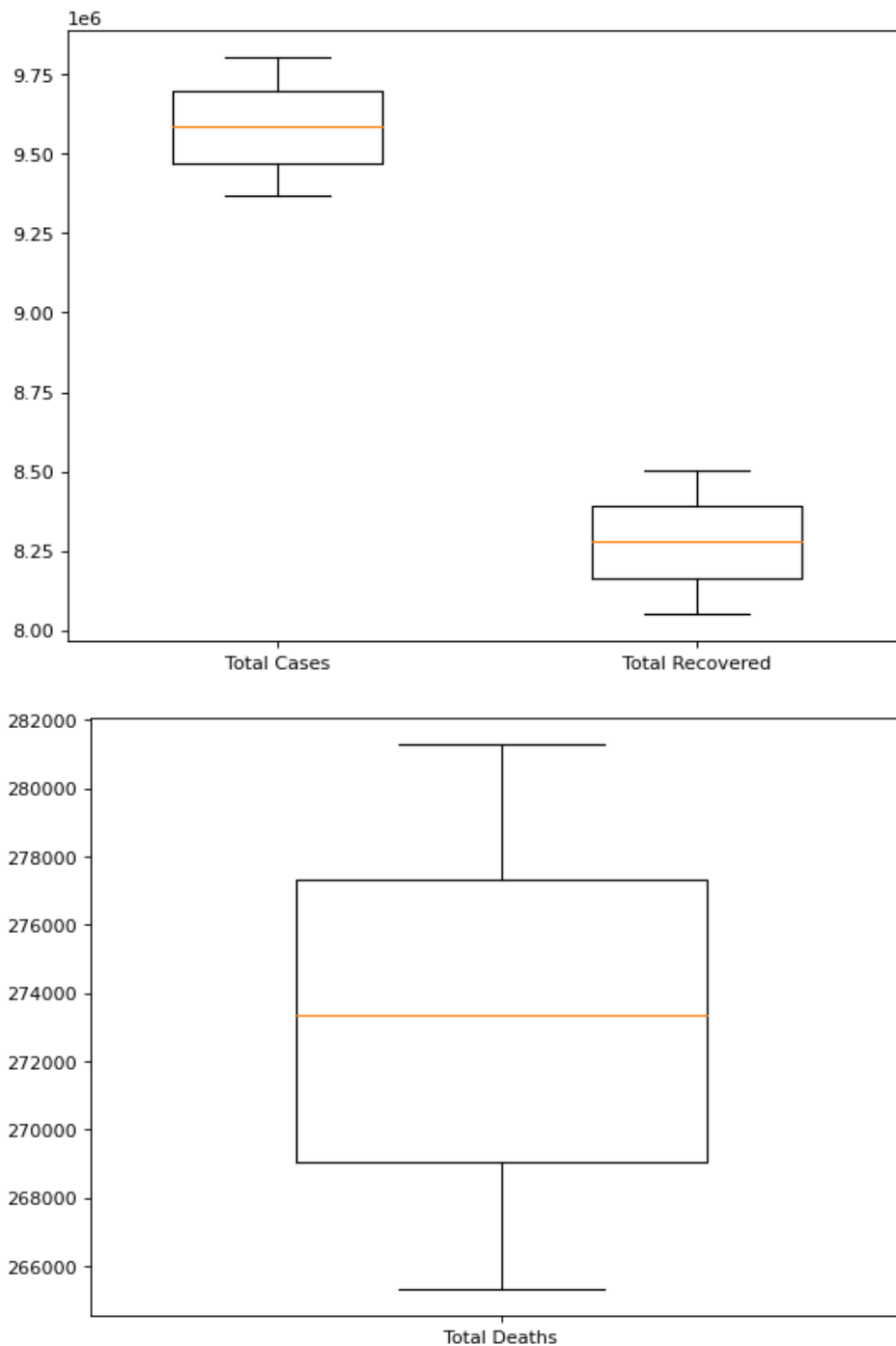
From the charts above, we can draw the conclusion that the distribution of the three attributes all revolve around the interval 0 and 500.000 cases. Little frequency in the middle, focused mainly on the left side of the chart.

There are more than 140 countries with at least 500,000 cases.





**Three boxplots** show some statistical measurements:
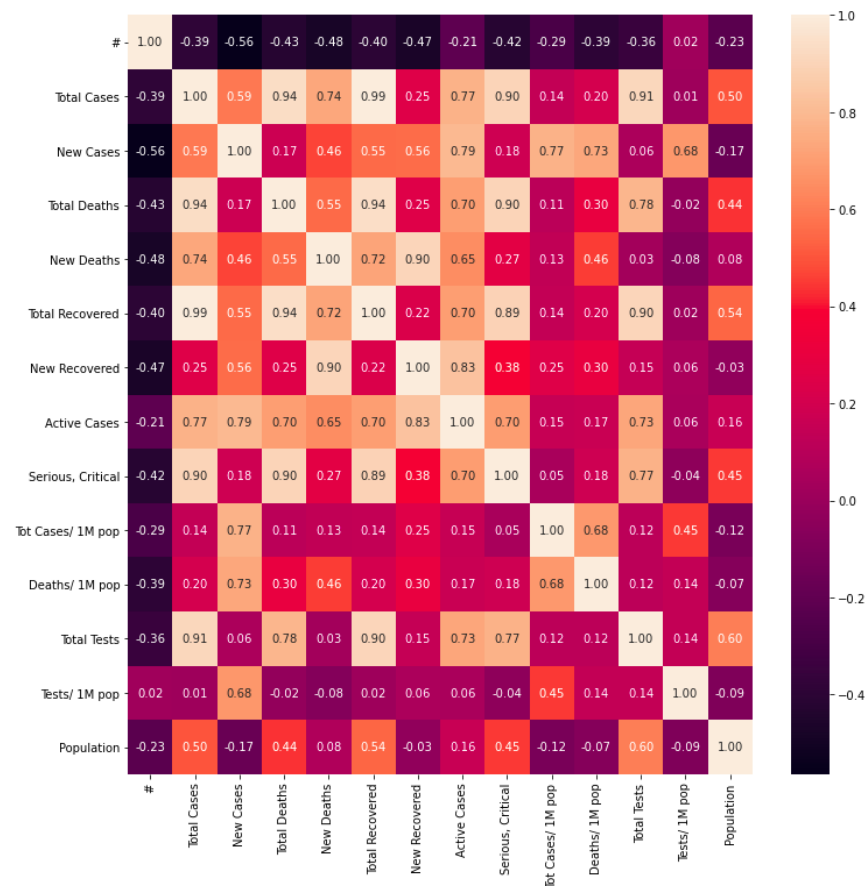- "Total Cases" has a median around 950.000 and 975.000 cases mean that

fifty-percentages of countries in that range. It shows that the whole world has a very high number of infections.

- "Total Recovered" has a median around 825.000 and 850.000 cases mean that fifty-percentages of countries in that range.
- "Total Deaths" has a median around 272.000 and 274.000 cases mean that fifty-percentages of countries in that range.

## e. Relation Plots

First, we use the **Correlation Matrix** to find correlational relationships of whole attributes.



We must ignore all "Total" prefix columns because "Total" prefix columns only increment if we consider it, it makes no sense.
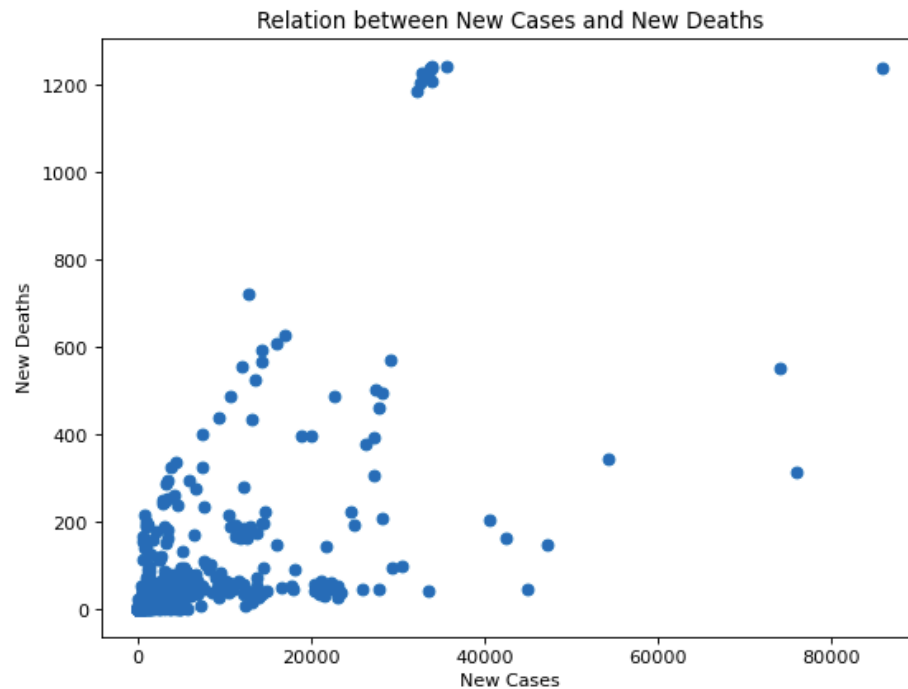
We can see that, "New Deaths" and "New Recovered" pair have the greatest correlational value of all. That means two columns have a relationship with each other.

In Relation Plots, we use **Scatter Charts** to find the relationship between two attributes of our data in order to find if the attributes covariate or inverse ?

In this example, we find the relationship of some pair attributes: "New Cases" and "New

Deaths", "New Deaths" and "New Recovered", "New Cases" and "New Recovered". Data is collected from all countries in the world from 22-11-2021 until 09-12-2021.
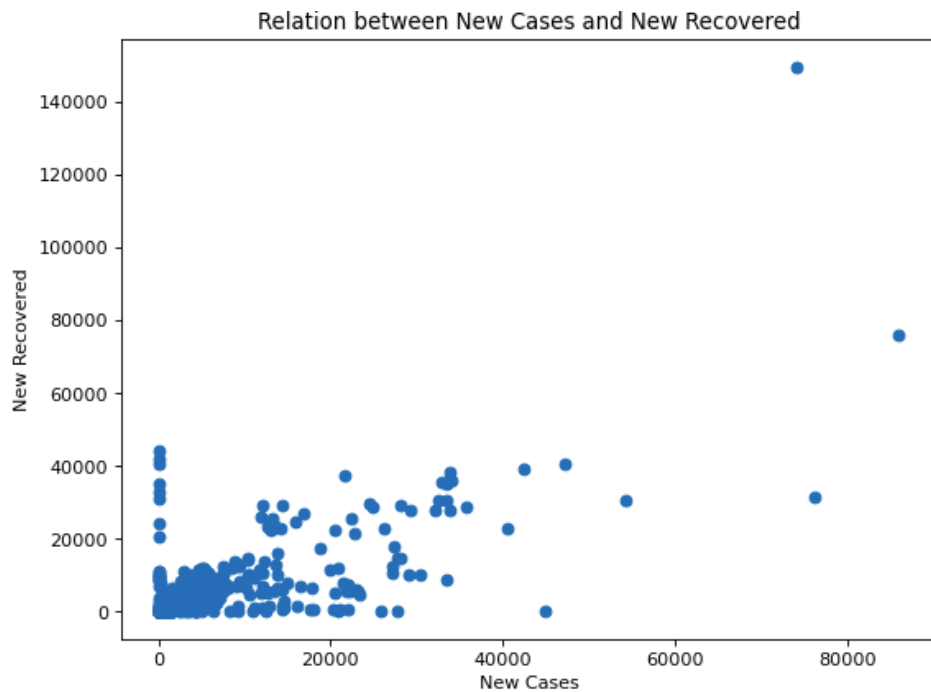


"New Cases" and "New Deaths" are not covariate. Some points of "New Deaths" increase but not "New Cases" and vice versa.
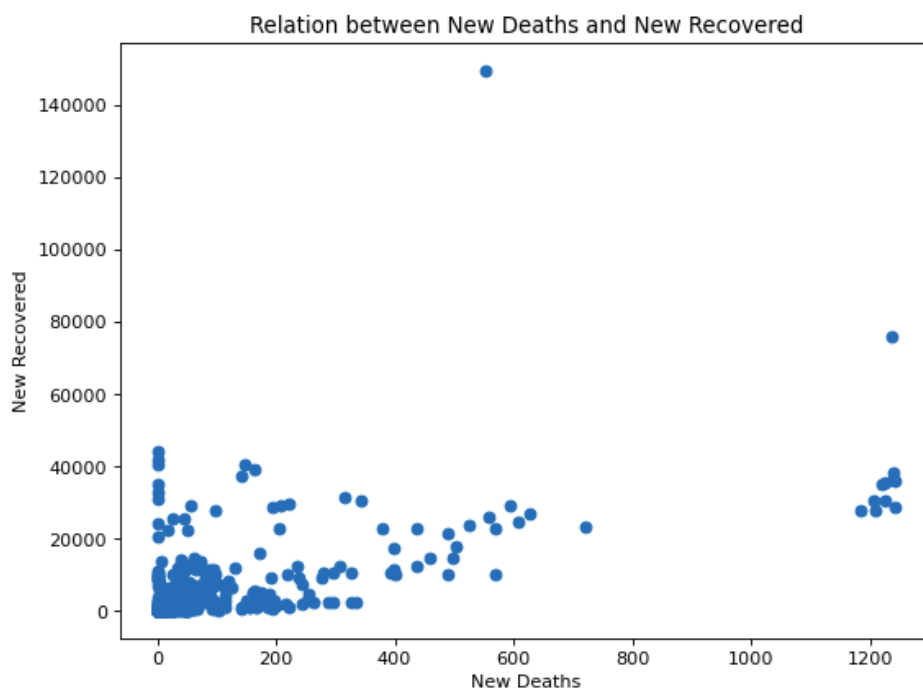
The number of points is most concentrated in the left region, showing countries with the number of deaths per day about 200 cases and the number of new cases per day about 20000 cases.

Relation between New Cases and New Recovered

Looking at the above chart, we can easily see that "New Cases" and "New Recovered" may be covariates. Because when "New Recovered" increases, so do "New Cases". Data is plotted, maybe lying on a straight line.



Relation between New Deaths and New Recovered

And the same as "New Deaths" and "New Recovered".

The number of points is most concentrated in the left region, showing countries with the number of recovered cases per day about 20000 cases and the number of death cases per day about 200 cases. And there are some countries with very high death cases around 1200 a day.

# 4.    References

https://airflow.apache.org/
https://www.applydatascience.com/airflow/airflow-tutorial-introduction/
https://devdocs.io/matplotlib~3.1/
Topic 04 - Visualization - Lecture document