# Project Documentation: Covid-19 Data Analysis

Ignatius S. Condro Atmawan

# Introduction

In the Covid-19 Data Analysis, the data procedure is like this:
1. Select variables of interest, from the data
2. Clean the data
3. Calculate descriptive statistics
4. Visualize the data with diagrams or other images

The data source comes from
https://www.kaggle.com/datasets/hendratno/covid19-indonesia/data
However, the data has noises after it is downloaded.

In order to process the data, the class CovidData is developed. This class is put in the script "covid_data_process.py" and the main file "covid19_main.py" is using the methods from the class to process Covid-19 data.

# CovidData Class

## Attributes

The attributes for this class are
- file_name:str, as the name of the file to be analyzed
- current_directory: str,
- output_directory: str
- df: pandas data frame
- columns: list, containing the list of columns of the data frame
- df_oi_dict: dict, containing the dictionary of the data frame of interest
- df_oi_statistics: dict, containing the dictionary of the descriptive statistics

## Methods

Methods for selecting variables
- select_column(self, list_of_columns: list)
- get_dataframe_of_interest_based_on_string()

Methods for cleaning data
- omit_empty_data(self)
- clean_numeric_data(self, list_of_columns: list)

Methods for calculating statistics
- calculate_descriptive_statistics(self, list_of_columns: list)
- save_descriptive_statistics(self)

Methods for data visualization
- plot_from_saved_dict(self, x_col: str, y_col: str, dict_key: str)

Other methods
- read_csv_data(self, file_name: str)
- save_csv_data(self, file_name: str)
- assign_directory(self)

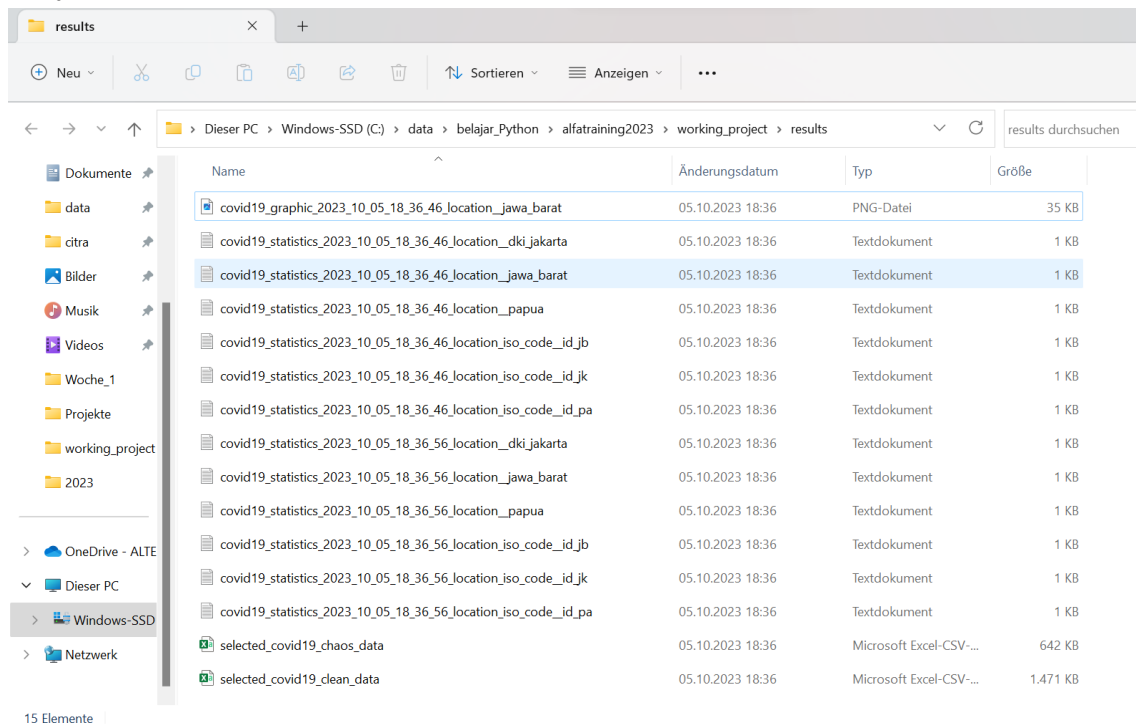# Data Processing

The data processing procedures is
1. Open two CSV files, with clean data as control and chaos data to test the data filtering/cleaning functions
   - "covid_19_indonesia_time_series_all.csv"
   - "Ignatius_covid_19_indonesia_time_series_all_chaos.csv"
2. Select the variables of interest:
   - for clean data: "New Cases", "New Deaths", "New Recovered", "New Active Cases"
   - for chaos data: "New Cases.1", "New Deaths.1", "New Recovered.1", "New Active Cases.1"
   - for both data: "Date", "Location ISO Code", "Location",
3. Cleaning the data using "clean_numeric_data" and "omit_empty_data"
4. Select the location of interest, e.g. "Jawa Barat", "DKI Jakarta", "Papua"
5. Calculate statistics based on our selection
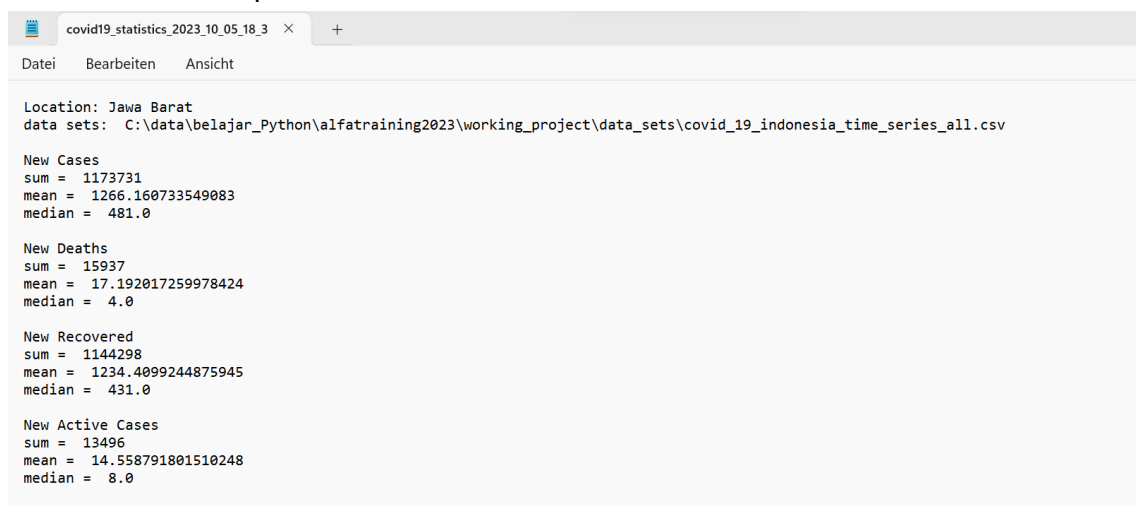6. Visualize the results

# Results

The results will be saved with date and time:
- The descriptive statistics are saved as text files
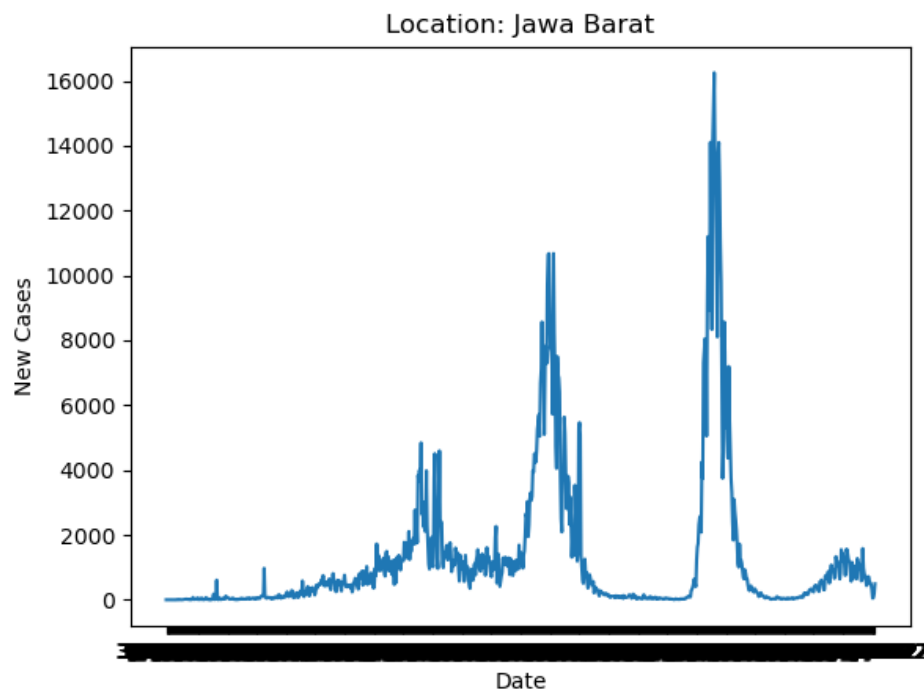- The data visualization is saved in figures of PNG format,

They are like this



The result of descriptive statistics is like this



The visualization result is like this

Location: Jawa Barat

The Pycharm logs will look like this



# Discussion

The data processing can clean the data, calculate statistics, and visualize the data. However the data cleaning process only deals with numerical data (integer) and it throws away a lot of data samples. If the string data can be cleaned, more data may be saved and not thrown away. The data visualization only shows minimum requirements, but the axis are not clearly visible.