

...

Soren DeOrlow
IDSN 599, Spring 2021
Deorlow@usc.edu
Final Project Part 2
...

PART 1

Domain: What is the type/nature of the problem that you are thinking to solve

Through the process of searching for data sources, there has been an abundance of data surrounding covid-19 along with interesting projects attempting to draw new insights from this vast data. I seek to become aware of the work that is happening in this space and through prototyping, gain new insight on data patterns around covid-19.

Dataset: Where did you get your dataset? How big is it – how many rows and attributes? Why/how does it contain the data you need to solve the problem you are thinking of solving?

I have found numerous datasets on Kaggle with fascinating data on covid 19, including lung x-ray imaging and covid vaccine related tweets. The cleanest dataset that I found maps vaccine data globally in relation to total population as well as covid fatalities. It contains 21k rows of data on global patterns. I also discovered that Roche has launched a project called the uncover Covid-19 challenge focused on data exploration and research. This project has launched 18 sub tasks inviting data research contributions to provide new perspective and answers to questions such as, "Predicting illness severity in a particular patient or demographic." I have chosen two similar datasets that I am working with to gain greater insight into the patterns of covid-19 from a public health standpoint and how various populations around the world have responded to vaccines and other measures.

Problem Type: Are you creating a predictor, a classifier, or something else? Why do you think this is the way to go to solve the problem?

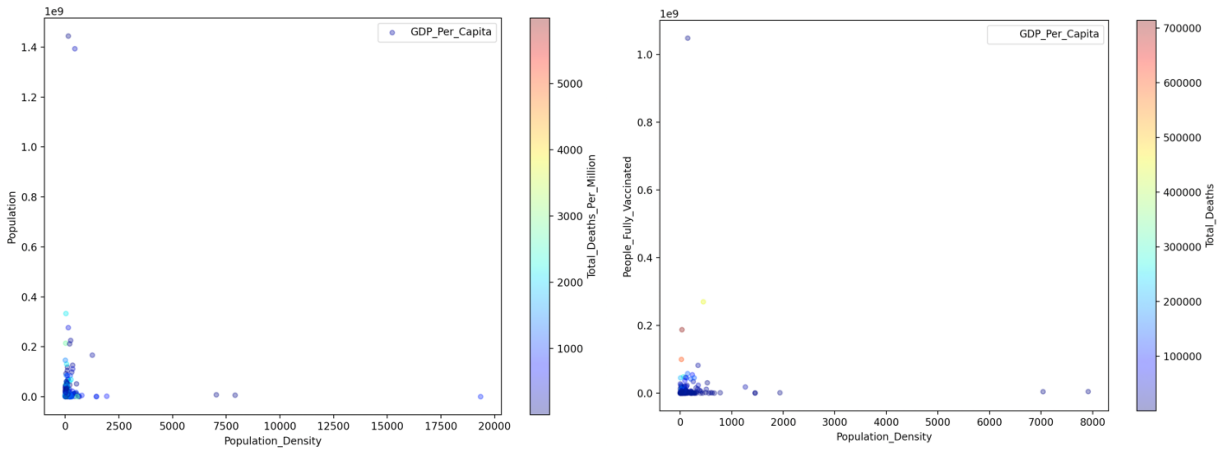
I plan to use classification to gain insight into how different parts of the world are approaching vaccination and the velocity of vaccination over time. I also am looking at the relationships of GDP, vaccination and covid morbidity. I'm less concerned with the relationship between covid morbidity, but I will not rule out any insights that might be found. Ultimately, I would like to see changes over time.

Attributes: What are the attributes of the dataset? Which are numeric and which are text? Do you have any missing values for attributes?

As I was exploring various datasets, I was careful to select datasets that were robust and consisted primarily of numeric data, with limited missing attributes. The datasets that I have chosen, list countries and how they are doing in the fight against covid-19. Below is a sample of both dataset including a complete view of the column headers.

Dataset 1

| | country | iso_code | date | total_vaccinations | people_vaccinated | people_fully_vaccinated | New_deaths | population | ratio |
|---|-------------|----------|---------|--------------------|-------------------|-------------------------|------------|------------|-------------|
| 0 | Afghanistan | AFG | 5/11/21 | 504502 | 448878 | 55624 | 12 | 40146987 | 1.118086396 |
| 1 | Afghanistan | AFG | 5/20/21 | 547901 | 470341 | 77560 | 10 | 40146987 | 1.171547444 |
| 2 | Afghanistan | AFG | 5/24/21 | 573277 | 476367 | 96910 | 10 | 40146987 | 1.186557288 |
| 3 | Afghanistan | AFG | 5/26/21 | 590454 | 479372 | 111082 | 19 | 40146987 | 1.194042283 |
| 4 | Afghanistan | AFG | 5/27/21 | 593313 | 479574 | 113739 | 14 | 40146987 | 1.194545434 |
| 5 | Afghanistan | AFG | 5/30/21 | 600152 | 480226 | 119926 | 20 | 40146987 | 1.196169466 |



PART 2

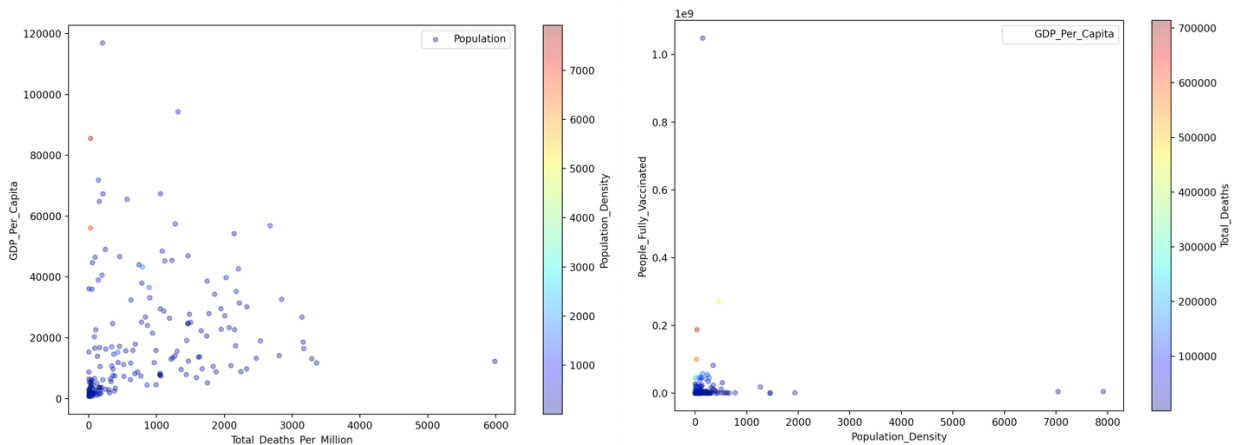
ML Algorithm: You are to choose three ML algorithms and describe why you have selected those algorithms. Justify your choices of ML algorithms.

I will be using `pd.get_dummies` (one hot encoder) to categorize countries

I will be using `cartopy.io.shapereader` in order to populate a global map with my dataset.

I will be using `matplotlib.pyplot` to evaluate and cluster my dataset to reveal relevant patterns.

Correlations: Show some numbers, or screen shots of plots to show which pairs of your data attribute correlations are most likely to be of value in your analysis. I suggest also showing plots as it really shows the quality of the correlation, instead of just numeric values.



Within the first dataset there is a correlation between GDP Per Capita and People Fully Vaccinated.

Transformers: Discuss each of your transformers, whether custom or not. Explain why you chose this transformer and what are each of them doing. I suggest using at least the `Pipeline` class. You may also want to use the `ColumnTransformer` class if you have categorical (text) data.

I will be using `pd.get_dummies` (one hot encoder) to transform nominal categorical data into categorized countries. This transformation will allow us to involve countries as a data source.

For my second dataset, I will create a custom transformer that will convert the categorical names of the different types of variants into a data source.

General notes.

The two datasets that I have selected show great promise. I am excited about exploring the migration of variants across the globe and the conditions that have enabled a variant to proliferate.

AutoSave OFF covid-variants

Home Insert Draw Page Layout Formulas Data Review Tell me Share Comments

Paste Font Alignment Number Conditional Formatting Format as Table Cells Editing Analyze Data Create and Share Adobe PDF

Possible Data Loss Some features might be lost if you save this workbook in the comma-delimited (.csv) format... Save As...

H15 fx

| | A | B | C | D | E | F |
|----|----------|----------|-----------|---------------|----------------|---------------------|
| 1 | location | date | variant | num_sequences | perc_sequences | num_sequences_total |
| 2 | Angola | 12/21/20 | B.1.160 | 0 | 0 | 93 |
| 3 | Angola | 12/21/20 | B.1.620 | 0 | 0 | 93 |
| 4 | Angola | 12/21/20 | B.1.258 | 0 | 0 | 93 |
| 5 | Angola | 12/21/20 | B.1.221 | | | 93 |
| 6 | Angola | 12/21/20 | B.1.1.302 | | | 93 |
| 7 | Angola | 12/21/20 | B.1.1.277 | | | 93 |
| 8 | Angola | 12/21/20 | B.1.1.519 | | | 93 |
| 9 | Angola | 12/21/20 | B.1.367 | | | 93 |
| 10 | Angola | 12/21/20 | B.1.177 | 0 | 0 | 93 |
| 11 | Angola | 12/21/20 | Beta | 68 | 73.12 | 93 |
| 12 | Angola | 12/21/20 | Alpha | 0 | 0 | 93 |
| 13 | Angola | 12/21/20 | Gamma | 0 | 0 | 93 |
| 14 | Angola | 12/21/20 | Delta | 0 | 0 | 93 |
| 15 | Angola | 12/21/20 | Kappa | 0 | 0 | 93 |
| 16 | Angola | 12/21/20 | Epsilon | | | 93 |
| 17 | Angola | 12/21/20 | Eta | 1 | 1.08 | 93 |
| 18 | Angola | 12/21/20 | Iota | | | 93 |
| 19 | Angola | 12/21/20 | Lambda | | | 93 |
| 20 | Angola | 12/21/20 | Mu | | | 93 |

Ready 110%

In [12]:

```

import pandas as pd
import os
import matplotlib.pyplot as plt

COVID19_PATH = os.path.join("COVID-19_Global_Dataset.csv")

def load_covid19_data(covid19_path = COVID19_PATH):
    return pd.read_csv(csv_path)

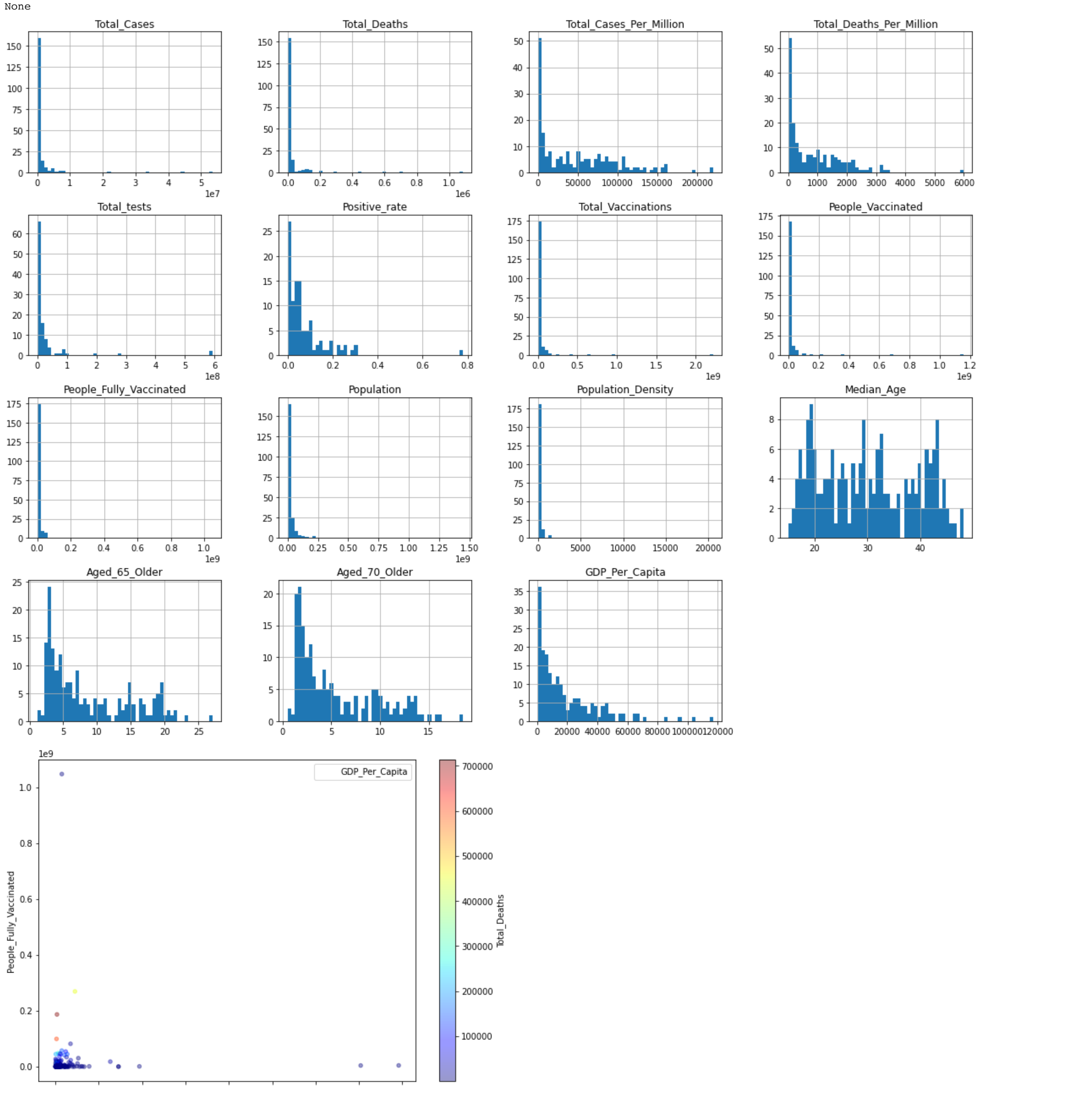
covid19 = pd.read_csv(COVID19_PATH)
covid19.info() #The info method is usefull to get a quick description of the data
print(covid19.info())
#print(covid19["people_vaccinated"].value_counts()) #Shows what categories exist and how many districts belong to each category
#print(covid19.describe()) #This method shows a summary of the numerical attributes
covid19.hist(bins=50, figsize=(20,15)) #shows the number of instances (vertical axis) that have a given value range
plt.show() #Plots a histogram for each numerical attribute
#covid19.plot(kind="scatter", x="GDP_Per_Capita", y="Total_Deaths", alpha=0.1)
#covid19.plot(kind="scatter", x="Total_Cases_Per_Million", y="GDP_Per_Capita", alpha=0.1)
covid19.plot(kind="scatter", x="Population_Density", y="People_Fully_Vaccinated", alpha=0.4, label="GDP_Per_Capita", figsize=(10,7),c="Total_Deaths", cmap=plt.get_cmap("magma"))
#print(covid19.median().values)
plt.legend()
plt.show() #This is required in PyCharm - not shown in book

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 216 entries, 0 to 215
Data columns (total 18 columns):
# Column Non-Null Count Dtype
---
0 Continent 214 non-null object
1 Country 216 non-null object
2 Last_Updated_Date 216 non-null object
3 Total_Cases 195 non-null float64
4 Total_Deaths 188 non-null float64
5 Total_Cases_Per_Million 195 non-null float64
6 Total_Deaths_Per_Million 188 non-null float64
7 Total_tests 104 non-null float64
8 Positive_rate 105 non-null float64
9 Total_Vaccinations 199 non-null float64
10 People_Vaccinated 194 non-null float64
11 People_Fully_Vaccinated 196 non-null float64
12 Population 216 non-null int64
13 Population_Density 204 non-null float64
14 Median_Age 189 non-null float64
15 Aged_65_Older 187 non-null float64
16 Aged_70_Older 188 non-null float64
17 GDP_Per_Capita 191 non-null float64
dtypes: float64(14), int64(1), object(3)
memory usage: 30.5+ KB
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 216 entries, 0 to 215
Data columns (total 18 columns):
# Column Non-Null Count Dtype
---
0 Continent 214 non-null object
1 Country 216 non-null object
2 Last_Updated_Date 216 non-null object
3 Total_Cases 195 non-null float64
4 Total_Deaths 188 non-null float64
5 Total_Cases_Per_Million 195 non-null float64
6 Total_Deaths_Per_Million 188 non-null float64
7 Total_tests 104 non-null float64
8 Positive_rate 105 non-null float64
9 Total_Vaccinations 199 non-null float64
10 People_Vaccinated 194 non-null float64
11 People_Fully_Vaccinated 196 non-null float64
12 Population 216 non-null int64
13 Population_Density 204 non-null float64
14 Median_Age 189 non-null float64
15 Aged_65_Older 187 non-null float64
16 Aged_70_Older 188 non-null float64
17 GDP_Per_Capita 191 non-null float64
dtypes: float64(14), int64(1), object(3)
memory usage: 30.5+ KB
None

```



In [11]:

```

import pandas as pd
import os
import matplotlib.pyplot as plt

COVID19_PATH = os.path.join("COVID-19_Global_Dataset.csv")

def load_covid19_data(covid19_path = COVID19_PATH):
    return pd.read_csv(csv_path)

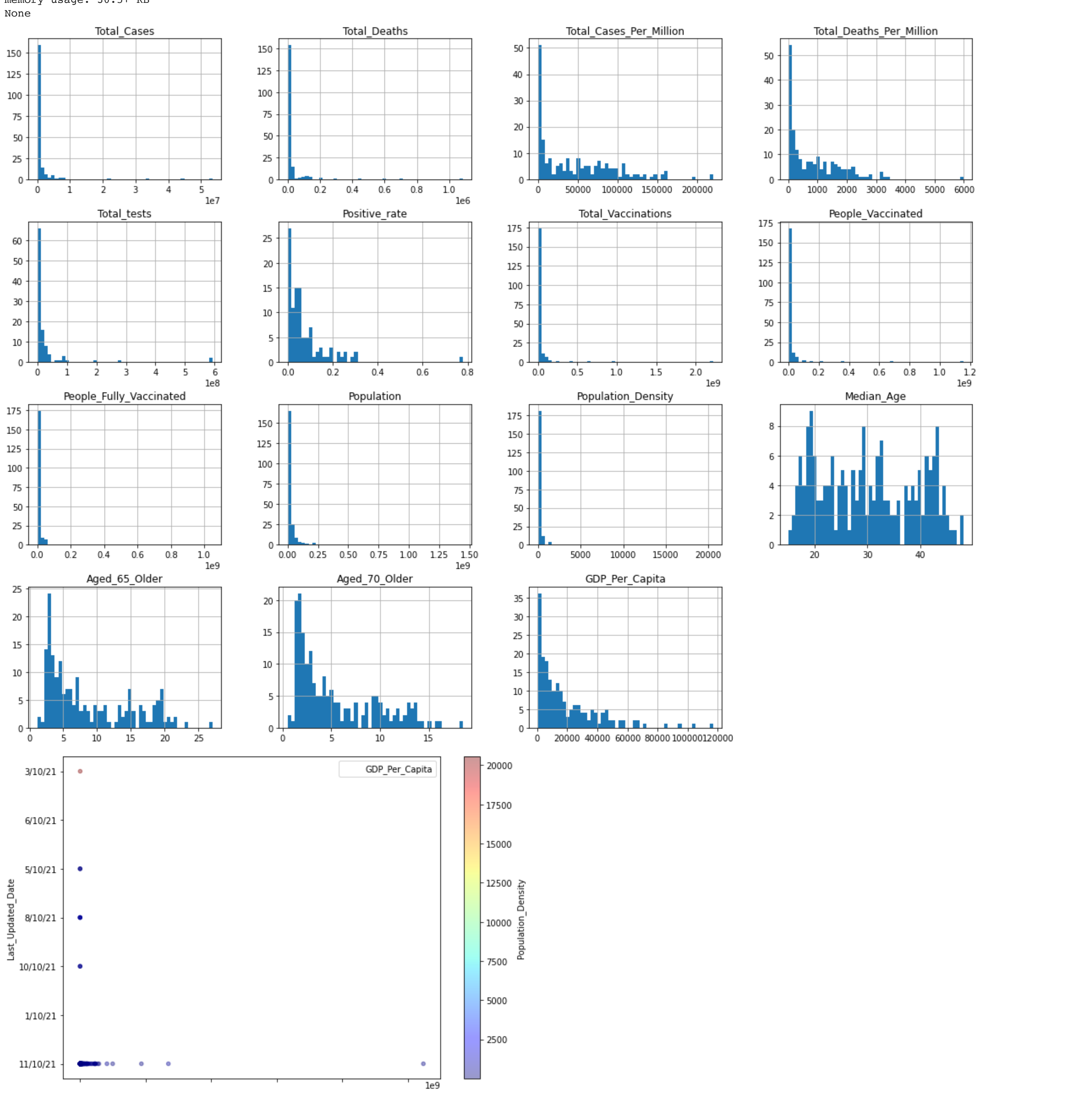
covid19 = pd.read_csv(COVID19_PATH)
covid19.info() #The info method is usefull to get a quick description of the data
print(covid19.info())
#print(covid19["people_vaccinated"].value_counts()) #Shows what categories exist and how many districts belong to each category
#print(covid19.describe()) #This method shows a summary of the numerical attributes
covid19.hist(bins=50, figsize=(20,15)) #shows the number of instances (vertical axis) that have a given value range
plt.show() #Plots a histogram for each numerical attribute
#covid19.plot(kind="scatter", x="GDP_Per_Capita", y="Total_Deaths", alpha=0.1)
#covid19.plot(kind="scatter", x="Total_Cases_Per_Million", y="GDP_Per_Capita", alpha=0.1)
covid19.plot(kind="scatter", x="People_Fully_Vaccinated", y="Last_Updated_Date", alpha=0.4, label="GDP_Per_Capita", figsize=(10,7),c="Population_Density", cmap=plt.get_cmap("magma"))
#print(covid19.median().values)
plt.legend()
plt.show() #This is required in PyCharm - not shown in book

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 216 entries, 0 to 215
Data columns (total 18 columns):
# Column Non-Null Count Dtype
---
0 Continent 214 non-null object
1 Country 216 non-null object
2 Last_Updated_Date 216 non-null object
3 Total_Cases 195 non-null float64
4 Total_Deaths 188 non-null float64
5 Total_Cases_Per_Million 195 non-null float64
6 Total_Deaths_Per_Million 188 non-null float64
7 Total_tests 104 non-null float64
8 Positive_rate 105 non-null float64
9 Total_Vaccinations 199 non-null float64
10 People_Vaccinated 194 non-null float64
11 People_Fully_Vaccinated 196 non-null float64
12 Population 216 non-null int64
13 Population_Density 204 non-null float64
14 Median_Age 189 non-null float64
15 Aged_65_Older 187 non-null float64
16 Aged_70_Older 188 non-null float64
17 GDP_Per_Capita 191 non-null float64
dtypes: float64(14), int64(1), object(3)
memory usage: 30.5+ KB
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 216 entries, 0 to 215
Data columns (total 18 columns):
# Column Non-Null Count Dtype
---
0 Continent 214 non-null object
1 Country 216 non-null object
2 Last_Updated_Date 216 non-null object
3 Total_Cases 195 non-null float64
4 Total_Deaths 188 non-null float64
5 Total_Cases_Per_Million 195 non-null float64
6 Total_Deaths_Per_Million 188 non-null float64
7 Total_tests 104 non-null float64
8 Positive_rate 105 non-null float64
9 Total_Vaccinations 199 non-null float64
10 People_Vaccinated 194 non-null float64
11 People_Fully_Vaccinated 196 non-null float64
12 Population 216 non-null int64
13 Population_Density 204 non-null float64
14 Median_Age 189 non-null float64
15 Aged_65_Older 187 non-null float64
16 Aged_70_Older 188 non-null float64
17 GDP_Per_Capita 191 non-null float64
dtypes: float64(14), int64(1), object(3)
memory usage: 30.5+ KB
None

```



In []: