

Covid - 19 Data Analysis Project using Python

```
In [2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import matplotlib.pyplot as plt
import seaborn as sns
import matplotlib.pyplot as plt

In [7]: url="https://raw.githubusercontent.com/SE5688/Datasets/main/covid-data.csv"
df=pd.read_csv(url)

Out[7]:
```

iso_code	continent	location	date	total_cases	new_cases	new_cases_smoothed	total_deaths	new_deaths	new_deaths_smoothed	...	gdp_per_capita	extreme_poverty	cardiovasc_death_rate	diabetes
0	AFG	Asia	Alghanistan	31/12/19	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN	597.029
1	AFG	Asia	Alghanistan	01/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN	597.029
2	AFG	Asia	Alghanistan	03/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN	597.029
3	AFG	Asia	Alghanistan	04/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN	597.029
...
57389	NaN	NaN	International	13/11/20	696.0	NaN	NaN	7.0	NaN	NaN	...	NaN	NaN	NaN
57390	NaN	NaN	International	14/11/20	696.0	NaN	NaN	7.0	NaN	NaN	...	NaN	NaN	NaN
57391	NaN	NaN	International	15/11/20	696.0	NaN	NaN	7.0	NaN	NaN	...	NaN	NaN	NaN
57392	NaN	NaN	International	16/11/20	696.0	NaN	NaN	7.0	NaN	NaN	...	NaN	NaN	NaN
57393	NaN	NaN	International	17/11/20	696.0	NaN	NaN	7.0	NaN	NaN	...	NaN	NaN	NaN

57394 rows × 15 columns

High Level Data Understanding:

```
In [8]: #Find no. of rows & columns in the dataset

print("no_of_rows = ",df.shape[0])
print("no_of_cols = ",df.shape[1])

no_of_rows = 57394
no_of_cols = 49

In [9]: #Data types of columns

df.dtypes

Out[9]:
```

iso_code	continent	location	date	total_cases	new_cases	new_cases_smoothed	total_deaths	new_deaths	new_deaths_smoothed	...	gdp_per_capita	extreme_poverty	cardiovasc_death_rate	diabetes
0	AFG	Asia	Alghanistan	31/12/19	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN	597.029
1	AFG	Asia	Alghanistan	01/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN	597.029
2	AFG	Asia	Alghanistan	03/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN	597.029
3	AFG	Asia	Alghanistan	04/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN	597.029
...
57389	NaN	NaN	International	13/11/20	696.0	NaN	NaN	7.0	NaN	NaN	...	NaN	NaN	NaN
57390	NaN	NaN	International	14/11/20	696.0	NaN	NaN	7.0	NaN	NaN	...	NaN	NaN	NaN
57391	NaN	NaN	International	15/11/20	696.0	NaN	NaN	7.0	NaN	NaN	...	NaN	NaN	NaN
57392	NaN	NaN	International	16/11/20	696.0	NaN	NaN	7.0	NaN	NaN	...	NaN	NaN	NaN
57393	NaN	NaN	International	17/11/20	696.0	NaN	NaN	7.0	NaN	NaN	...	NaN	NaN	NaN

57394 rows × 15 columns

```
In [10]: #Info of data in dataframe

df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 57394 entries, 0 to 57393
Data columns (total 49 columns):
 # Column Non-Null Count Dtype
---
 0 iso_code 57871 non-null object
 1 continent 56748 non-null object
 2 location 57394 non-null object
 3 date 57394 non-null object
 4 total_cases 53758 non-null float64
 5 new_cases 56465 non-null float64
 6 new_cases_smoothed 55652 non-null float64
 7 total_deaths 43568 non-null float64
 8 new_deaths 56652 non-null float64
 9 new_deaths_smoothed 55652 non-null float64
10 total_cases_per_million 53471 non-null float64
11 new_cases_per_million 56401 non-null float64
12 total_deaths_per_million 53471 non-null float64
13 new_deaths_per_million 56401 non-null float64
14 reproduction_rate 55887 non-null float64
15 icu_patients 56401 non-null float64
16 icu_patients_per_million 55887 non-null float64
17 hosp_patients 56401 non-null float64
18 hosp_patients_per_million 55887 non-null float64
19 weekly_icu_admissions 4480 non-null float64
20 weekly_hosp_admissions 4480 non-null float64
21 weekly_icu_admissions_per_million 5085 non-null float64
22 weekly_hosp_admissions_per_million 5085 non-null float64
23 new_tests 357 non-null float64
24 total_tests_per_thousand 645 non-null float64
25 new_tests_smoothed 645 non-null float64
26 new_tests_smoothed_per_thousand 645 non-null float64
27 tests_per_case 2287 non-null float64
28 positive_rate 21787 non-null float64
29 population 2287 non-null float64
30 population_density 21787 non-null float64
31 median_age 24632 non-null float64
32 aged_65_and_over 24632 non-null float64
33 aged_70_and_over 23211 non-null float64
34 population 47847 non-null float64
35 population_density 57871 non-null float64
36 median_age 54371 non-null float64
37 aged_65_and_over 51834 non-null float64
38 aged_70_and_over 50265 non-null float64
39 gdp_per_capita 50768 non-null float64
40 extreme_poverty 50367 non-null float64
41 cardiovasc_death_rate 2571 non-null float64
42 diabetes_prevalence 51813 non-null float64
43 female_smokers 39156 non-null float64
44 male_smokers 24176 non-null float64
45 handwashing_facilities 46836 non-null float64
46 hospital_beds_per_thousand 56336 non-null float64
47 life_expectancy 49247 non-null float64
48 human_development_index
49 dtype: object
memory usage: 21.5+ MB

In [11]: #describe of data in dataframe

df.describe()

Out[11]:
```

	total_cases	new_cases	new_cases_smoothed	total_deaths	new_deaths	new_deaths_smoothed	total_cases_per_million	new_cases_per_million	new_cases_smoothed_per_million	total_deaths_per_million
count	5.375800e+04	56465.000000	55652.000000	4.436500e+04	56465.000000	55652.000000	53471.000000	56401.000000	55587.000000	44096.000000
mean	1.677794e+05	1953.576941	1902.431953	6.958639e+03	47.054317	46.835439	3139.099992	39.344804	38.403211	95.335291
std	1.695038e+06	18289.650340	17777.391785	5.078061e+04	390.853776	378.272784	6183.455260	133.989155	101.648441	106.216999
min	1.000000e+00	4261.000000	0.50000000	1.000000e+00	1918.000000	232.143000	0.0010000	2212.540000	269.978000	0.0000000
25%	1.800000e+02	0.000000	0.00000000	1.000000e+01	0.000000	0.000000	90.033000	0.000000	0.260000	3.977750
50%	2.070000e+03	14.000000	13.4200000	8.400000e+01	0.000000	0.286000	593.663000	2.194000	3.863000	20.383000
75%	2.235675e+04	235.000000	245.286000	7.270000e+02	4.000000	4.000000	3487.050000	25.941000	29.678000	90.512150
max	5.515465e+07	646211.000000	584948.857000	1.328570e+06	10600.000000	9027.714000	76541.772000	8655.659000	2472.188000	1248.014000

8 rows × 11 columns

Low Level Data Understanding

```
In [12]: #Find count of unique values in location column.

df['location'].nunique()

Out[12]: 216

In [13]: #Find which continent has maximum frequency using values counts.

df['continent'].value_counts()

Out[13]:
```

continent	count
Europe	14828
Africa	13937
Asia	13528
North America	9116
South America	3484
Oceania	2235

name: continent, dtype: int64

```
In [14]: #Find maximum & mean value in 'total_cases'.

print("Maximum value: ",df['total_cases'].max())
print("Mean value: ",df['total_cases'].mean())

Maximum value: 55154651.0
Mean value: 167797.368875392

In [15]: #Find 25%,50% & 75% quartile value in 'total_deaths'.

print("25% ",df['total_deaths'].describe()[4])
print("50% ",df['total_deaths'].describe()[5])
print("75% ",df['total_deaths'].describe()[6])

25% 13.0
50% 84.0
75% 727.0

In [16]: #Find which continent has maximum 'human_development_index'.

df.groupby(['continent']).agg({'human_development_index':'max'})

Out[16]:
```

continent	human_development_index
Africa	0.797
Asia	0.933
Europe	0.963
North America	0.926
Oceania	0.939
South America	0.843

```
In [17]: #Find which continent has minimum 'gdp_per_capita'.

df.groupby(['continent']).agg({'gdp_per_capita':'min'})

Out[17]:
```

continent	gdp_per_capita
Africa	661.140
Asia	1479.147
Europe	5189.972
North America	1653.173
Oceania	2205.923
South America	6885.829

```
In [18]: # Filter the dataframe with only this columns ['continent','location','date','total_cases','total_deaths','gdp_per_capita','human_development_index']
# and update the data frame.

df=df[['continent','location','date','total_cases','total_deaths','gdp_per_capita','human_development_index']]
df

Out[18]:
```

continent	location	date	total_cases	total_deaths	gdp_per_capita	human_development_index	
0	Asia	Alghanistan	31/12/19	NaN	NaN	1803.987	0.498
1	Asia	Alghanistan	01/01/20	NaN	NaN	1803.987	0.498
2	Asia	Alghanistan	03/01/20	NaN	NaN	1803.987	0.498
3	Asia	Alghanistan	03/01/20	NaN	NaN	1803.987	0.498
4	Asia	Alghanistan	04/01/20	NaN	NaN	1803.987	0.498
...
57389	NaN	International	13/11/20	696.0	7.0	NaN	NaN
57390	NaN	International	14/11/20	696.0	7.0	NaN	NaN
57391	NaN	International	15/11/20	696.0	7.0	NaN	NaN
57392	NaN	International	16/11/20	696.0	7.0	NaN	NaN
57393	NaN	International	17/11/20	696.0	7.0	NaN	NaN

57394 rows × 8 columns

Data Cleaning

```
In [19]: #Remove all duplicates observations

df.drop_duplicates()

Out[19]:
```

continent	location	date	total_cases	total_deaths	gdp_per_capita	human_development_index	
0	Asia	Alghanistan	31/12/19	NaN	NaN	1803.987	0.498
1	Asia	Alghanistan	01/01/20	NaN	NaN	1803.987	0.498
2	Asia	Alghanistan	03/01/20	NaN	NaN	1803.987	0.498
3	Asia	Alghanistan	03/01/20	NaN	NaN	1803.987	0.498
4	Asia	Alghanistan	04/01/20	NaN	NaN	1803.987	0.498
...
57389	NaN	International	13/11/20	696.0	7.0	NaN	NaN
57390	NaN	International	14/11/20	696.0	7.0	NaN	NaN
57391	NaN	International	15/11/20	696.0	7.0	NaN	NaN
57392	NaN	International	16/11/20	696.0	7.0	NaN	NaN
57393	NaN	International	17/11/20	696.0	7.0	NaN	NaN

57394 rows × 8 columns

```
In [20]: #Find missing values in all columns

df.isnull()

Out[20]:
```

continent	location	date	total_cases	total_deaths	gdp_per_capita	human_development_index	
0	False	False	False	True	True	False	False
1	False	False	False	True	True	False	False
2	False	False	False	True	True	False	False
3	False	False	False	True	True	False	False
4	False	False	False	True	True	False	False
...
57389	True	False	False	False	False	True	True
57390	True	False	False	False	False	True	True
57391	True	False	False	False	False	True	True
57392	True	False	False	False	False	True	True
57393	True	False	False	False	False	True	True

57394 rows × 8 columns

```
In [21]: # Remove all observations where continent value is missing

df.dropna(subset=['continent'])

Out[21]:
```

continent	location	date	total_cases	total_deaths	gdp_per_capita	human_development_index	
0	Asia	Alghanistan	31/12/19	NaN	NaN	1803.987	0.498
1	Asia	Alghanistan	01/01/20	NaN	NaN	1803.987	0.498
2	Asia	Alghanistan	03/01/20	NaN	NaN	1803.987	0.498
3	Asia	Alghanistan	03/01/20	NaN	NaN	1803.987	0.498
4	Asia	Alghanistan	04/01/20	NaN	NaN	1803.987	0.498
...
56745	Africa	Zimbabwe	13/11/20	696.0	257.0	1899.775	0.535
56746	Africa	Zimbabwe	14/11/20	6785.0	257.0	1899.775	0.535
56745	Africa	Zimbabwe	15/11/20	6786.0	257.0	1899.775	0.535
56746	Africa	Zimbabwe	16/11/20	6786.0	257.0	1899.775	0.535
56747	Africa	Zimbabwe	17/11/20	8087.0	257.0	1899.775	0.535

56748 rows × 8 columns

```
In [22]: #Fill all missing values with 0

df.fillna(0)

Out[22]:
```

continent	location	date	total_cases	total_deaths	gdp_per_capita	human_development_index	
0	Asia	Alghanistan	31/12/19	0.0	0.0	1803.987	0.498
1	Asia	Alghanistan	01/01/20	0.0	0.0	1803.987	0.498
2	Asia	Alghanistan	03/01/20	0.0	0.0	1803.987	0.498
3	Asia	Alghanistan	03/01/20	0.0	0.0	1803.987	0.498
4	Asia	Alghanistan	04/01/20	0.0	0.0	1803.987	0.498
...
57389	0	International	13/11/20	696.0	7.0	0.000	0.000
57390	0	International	14/11/20	696.0	7.0	0.000	0.000
57391	0	International	15/11/20	696.0	7.0	0.000	0.000
57392	0	International	16/11/20	696.0	7.0	0.000	0.000
57393	0	International	17/11/20	696.0	7.0	0.000	0.000

57394 rows × 8 columns

```
In [ ]: #Convert date column in datetime format using pandas.to_datetime

df['date']=pd.to_datetime(df['date'])

In [26]: df

Out[26]:
```

continent	location	date	total_cases	total_deaths	gdp_per_capita	human_development_index	
0	Asia	Alghanistan	2019-12-31	NaN	NaN	1803.987	0.498
1	Asia	Alghanistan	2020-01-01	NaN	NaN	1803.987	0.498
2	Asia	Alghanistan	2020-02-01	NaN	NaN	1803.987	0.498
3	Asia	Alghanistan	2020-03-01	NaN	NaN	1803.987	0.498
4	Asia	Alghanistan	2020-04-01	NaN	NaN	1803.987	0.498
...
57389	NaN	International	2020-11-13	696.0	7.0	NaN	NaN
57390	NaN	International	2020-11-14	696.0	7.0	NaN	NaN
57391	NaN	International	2020-11-15	696.0	7.0	NaN	NaN
57392	NaN	International	2020-11-16	696.0	7.0	NaN	NaN
57393	NaN	International	2020-11-17	696.0	7.0	NaN	NaN

57394 rows × 8 columns

Data Aggregation:

```
In [36]: #Find max value in all columns using groupby function on 'continent' column

df.groupby(["continent"]).max().reset_index()

Out[36]:
```

continent	location	date	total_cases	total_deaths	gdp_per_capita	human_development_index	month	
0	Africa	Zimbabwe	2020-12-11	752269.0	20314.0	26382.287	0.797	12
1	Asia	Yemen	2020-12-11	8874290.0	130519.0	116935.600	0.933	12
2	Europe	Vatican	2020-12-11	1991233.0	52147.0	94277.965	0.963	12
3	North America	United States Virgin Islands	2020-12-11	11205486.0	247220.0	54225.446	0.926	12
4	Oceania	Wallis and Futuna	2020-12-11	27750.0	907.0	44648.710	0.939	12
5	South America	Venezuela	2020-12-11	5876464.0	166014.0	22767.037	0.843	12

```
In [31]: #Store the result in a new dataframe named 'df_groupby'.

df_groupby = df.groupby(["continent"]).max().reset_index()

In [32]: df_groupby

Out[32]:
```

continent	location	date	total_cases	total_deaths	gdp_per_capita	human_development_index	month	
0	Africa	Zimbabwe	2020-12-11	752269.0	20314.0	26382.287	0.797	12
1	Asia	Yemen	2020-12-11	8874290				