

```
In [ ]: import io

In [ ]: import pandas as pd

In [ ]:

In [ ]: import requests

In [ ]:

In [ ]: import numpy as np
```

1. import the dataset using pandas from url

```
In [9]: url = pd.read_csv("https://raw.githubusercontent.com/SR1608/Datasets/main/covid-data.csv")
url

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
0	AFG	Asia	Afghanistan	31/12/19	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
1	AFG	Asia	Afghanistan	01/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
2	AFG	Asia	Afghanistan	02/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
3	AFG	Asia	Afghanistan	03/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
4	AFG	Asia	Afghanistan	04/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
...
57389	NaN	NaN	International	13/11/20	696.0	NaN	NaN	NaN	7.0	NaN	NaN	NaN	NaN
57390	NaN	NaN	International	14/11/20	696.0	NaN	NaN	NaN	7.0	NaN	NaN	NaN	NaN
57391	NaN	NaN	International	15/11/20	696.0	NaN	NaN	NaN	7.0	NaN	NaN	NaN	NaN
57392	NaN	NaN	International	16/11/20	696.0	NaN	NaN	NaN	7.0	NaN	NaN	NaN	NaN
57393	NaN	NaN	International	17/11/20	696.0	NaN	NaN	NaN	7.0	NaN	NaN	NaN	NaN

57394 rows × 49 columns

2. high level data understanding

a. Find no. of rows & columns in the dataset b. Data types of columns. c. Info & describe of data in dataframe.

```
In [10]: url.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        44386 non-null float64
 8   new_deaths          56465 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 56465 non-null float64
12   new_cases_smoothed_per_million 55587 non-null float64
13   total_deaths_per_million 44096 non-null float64
14   new_deaths_per_million 56465 non-null float64
15   new_deaths_smoothed_per_million 55587 non-null float64
16   reproduction_rate    37696 non-null float64
17   icu_patients         4498 non-null float64
18   icu_patients_per_million 5695 non-null float64
19   hosp_patients_per_million 5895 non-null float64
20   weekly_icu_admissions 357 non-null float64
21   weekly_icu_admissions_per_million 357 non-null float64
22   weekly_icu_admissions_per_million 645 non-null float64
23   weekly_hosp_admissions_per_million 645 non-null float64
24   weekly_hosp_admissions_per_million 21787 non-null float64
25   total_tests          22917 non-null float64
26   new_tests            21787 non-null float64
27   total_tests_per_thousand 22917 non-null float64
28   new_tests_per_thousand 21787 non-null float64
29   new_tests_smoothed   24612 non-null float64
30   new_tests_smoothed_per_thousand 24612 non-null float64
31   tests_per_case       22892 non-null float64
32   positive_rate        23211 non-null float64
33   stringency_index     47847 non-null float64
34   population            56265 non-null float64
35   population_density   54371 non-null float64
36   median_age           51034 non-null float64
37   aged_65_and_over     33571 non-null float64
38   aged_70_and_over     58768 non-null float64
39   gdp_per_capita       58367 non-null float64
40   extreme_poverty      33571 non-null float64
41   cardiovascular_death_rate 51913 non-null float64
42   diabetes_prevalence  52881 non-null float64
43   female_smokers        39669 non-null float64
44   male_smokers          39156 non-null float64
45   handwashing_facilities 24176 non-null float64
46   hospital_beds_per_thousand 45926 non-null float64
47   life_expectancy      56336 non-null float64
48   human_development_index 49247 non-null float64
dtypes: float64(45), object(4)
memory usage: 21.5+ MB
```

3. low level data understanding

a. Find count of unique values in location column.

```
In [23]: url['location'].value_counts()

Out[23]:
United Kingdom    323
San Marino        323
Estonia           323
Greece            323
Brazil            323
Hong Kong        72
Solomon Islands  33
Wallis and Futuna 32
Marshall Islands 29
Vanuatu           7
Name: location, Length: 216, dtype: int64

b. Find which continent has maximum frequency using value counts.

In [14]: url['continent'].value_counts(dropna=False)

Out[14]:
Europe    14828
Africa    13637
Asia      13528
North America  9116
South America 3484
Oceania     2235
NaN         646
Name: continent, dtype: int64
Europe has the highest frequency of 14828

c. Find maximum & mean value in 'total_cases'.
```

```
In [26]: column_name="total_cases"
column_sum=url[column_name].sum()
a=column_sum/57393
print("The maximum and mean of total cases is: ",a)

The maximum and mean of total cases is:  157169.88857777882

e. Find which continent has maximum human development index

In [34]: url.groupby('continent')['human_development_index'].sum()

Out[34]:
continent    7222.777
Africa       9661.617
Asia         9152.325
North America 4494.381
Oceania       984.648
South America 2386.556
NaN          646.000
Name: human_development_index, dtype: float64
Continent with highest human development index is Europe that is 10817.325

f. Find which continent has minimum 'gdp_per_capita'.
```

```
In [35]: url.groupby('continent')['gdp_per_capita'].sum()

Out[35]:
continent    7.392128e+07
Africa       3.611561e+08
Asia         4.313814e+08
North America 1.536159e+08
South America 2.972206e+07
South America 4.378286e+07
Name: gdp_per_capita, dtype: float64
Continent with minimum gdp per capita is North America
```

4. 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.

```
In [37]: df.filter(items=['continent','location','date','total_cases','total_deaths','gdp_per_capita','human_development_index'])

Out[37]:
```

	continent	location	date	total_cases	total_deaths	gdp_per_capita	human_development_index
0	Asia	Afghanistan	31/12/19	NaN	NaN	1803.987	0.498
1	Asia	Afghanistan	01/01/20	NaN	NaN	1803.987	0.498
2	Asia	Afghanistan	02/01/20	NaN	NaN	1803.987	0.498
3	Asia	Afghanistan	03/01/20	NaN	NaN	1803.987	0.498
4	Asia	Afghanistan	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 × 7 columns

5. Data Cleaning

a. Remove all duplicates observations

```
In [6]: url.sort_values("date", inplace = True)
url.drop_duplicates(subset ="date",
                    keep = False, inplace = True)
url

Out[6]:
iso_code  continent  location  date  total_cases  new_cases  new_cases_smoothed  total_deaths  new_deaths  new_deaths_smoothed  ...  gdp_per_capita  extreme_poverty  cardiovascular_c
0 rows × 49 columns

we think there are no repeated observation in the dataframe

b. Find missing values in all columns

In [15]: import pandas as pd
import numpy as np
url = pd.read_csv("https://raw.githubusercontent.com/SR1608/Datasets/main/covid-data.csv")
url.fillna(method = 'pad')

Out[15]:
```

	iso_code	continent	location	date	total_cases	new_cases	new_cases_smoothed	total_deaths	new_deaths	new_deaths_smoothed	...	gdp_per_capita	extreme_poverty
0	AFG	Asia	Afghanistan	31/12/19	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
1	AFG	Asia	Afghanistan	01/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
2	AFG	Asia	Afghanistan	02/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
3	AFG	Asia	Afghanistan	03/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
4	AFG	Asia	Afghanistan	04/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
...
57389	OWID_WRL	Africa	International	13/11/20	696.0	-9.0	-1.286	7.0	1.0	0.143	...	15469.207	10.0
57390	OWID_WRL	Africa	International	14/11/20	696.0	-9.0	-1.286	7.0	1.0	0.143	...	15469.207	10.0
57391	OWID_WRL	Africa	International	15/11/20	696.0	-9.0	-1.286	7.0	1.0	0.143	...	15469.207	10.0
57392	OWID_WRL	Africa	International	16/11/20	696.0	-9.0	-1.286	7.0	1.0	0.143	...	15469.207	10.0
57393	OWID_WRL	Africa	International	17/11/20	696.0	-9.0	-1.286	7.0	1.0	0.143	...	15469.207	10.0

57394 rows × 49 columns

c. Remove all observations where continent column value is missing

```
In [17]: mod_url = url.dropna(how='any',
                           subset=['continent'])
print("Modified DataFrame : ")
display(mod_url)

Modified DataFrame :
```

	iso_code	continent	location	date	total_cases	new_cases	new_cases_smoothed	total_deaths	new_deaths	new_deaths_smoothed	...	gdp_per_capita	extreme_poverty
0	AFG	Asia	Afghanistan	31/12/19	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
1	AFG	Asia	Afghanistan	01/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
2	AFG	Asia	Afghanistan	02/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
3	AFG	Asia	Afghanistan	03/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
4	AFG	Asia	Afghanistan	04/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	1803.987	NaN
...
56743	ZWE	Africa	Zimbabwe	13/11/20	8696.0	29.0	36.000	255.0	0.0	1.000	...	1899.775	21.4
56744	ZWE	Africa	Zimbabwe	14/11/20	8765.0	69.0	42.000	257.0	2.0	1.000	...	1899.775	21.4
56745	ZWE	Africa	Zimbabwe	15/11/20	8795.0	21.0	41.143	257.0	0.0	0.857	...	1899.775	21.4
56746	ZWE	Africa	Zimbabwe	16/11/20	8795.0	0.0	36.429	257.0	0.0	0.571	...	1899.775	21.4
56747	ZWE	Africa	Zimbabwe	17/11/20	8897.0	111.0	48.000	257.0	0.0	0.429	...	1899.775	21.4

56748 rows × 49 columns

d. Fill all missing values with 0

```
In [19]: import pandas as pd
import numpy as np
url = pd.read_csv("https://raw.githubusercontent.com/SR1608/Datasets/main/covid-data.csv")
url.fillna(0)

Out[19]:
```

	iso_code	continent	location	date	total_cases	new_cases	new_cases_smoothed	total_deaths	new_deaths	new_deaths_smoothed	...	gdp_per_capita	extreme_poverty
0	AFG	Asia	Afghanistan	31/12/19	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0
1	AFG	Asia	Afghanistan	01/01/20	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0
2	AFG	Asia	Afghanistan	02/01/20	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0
3	AFG	Asia	Afghanistan	03/01/20	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0
4	AFG	Asia	Afghanistan	04/01/20	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0
...
57389	0	0	International	13/11/20	696.0	0.0	0.0	7.0	0.0	0.0	...	0.000	0.0
57390	0	0	International	14/11/20	696.0	0.0	0.0	7.0	0.0	0.0	...	0.000	0.0
57391	0	0	International	15/11/20	696.0	0.0	0.0	7.0	0.0	0.0	...	0.000	0.0
57392	0	0	International	16/11/20	696.0	0.0	0.0	7.0	0.0	0.0	...	0.000	0.0
57393	0	0	International	17/11/20	696.0	0.0	0.0	7.0	0.0	0.0	...	0.000	0.0

57394 rows × 49 columns

6. Date time format :

a. Convert date column in datetime format using pandas.to_datetime

```
In [6]: import pandas as pd
url = pd.read_csv("https://raw.githubusercontent.com/SR1608/Datasets/main/covid-data.csv")
a=url['date'] = pd.to_datetime(url['date'])
display(a)

0      2019-12-31
1      2020-01-01
2      2020-02-01
3      2020-03-01
4      2020-04-01
...
57389  2020-11-13
57390  2020-11-14
57391  2020-11-15
57392  2020-11-16
57393  2020-11-17
Name: date, Length: 57394, dtype: datetime64[ns]

b. Create new column month after extracting date from date column.
```

```
In [10]: import pandas as pd
url = pd.read_csv("https://raw.githubusercontent.com/SR1608/Datasets/main/covid-data.csv")
url['month'] = pd.to_datetime(url['date']).month
url

Out[10]:
```

	iso_code	continent	location	date	total_cases	new_cases	new_cases_smoothed	total_deaths	new_deaths	new_deaths_smoothed	...	extreme_poverty	cardiovasc_death
0	AFG	Asia	Afghanistan	31/12/19	NaN	0.0	NaN	NaN	0.0	NaN	...	NaN	59
1	AFG	Asia	Afghanistan	01/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	NaN	59
2	AFG	Asia	Afghanistan	02/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	NaN	59
3	AFG	Asia	Afghanistan	03/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	NaN	59
4	AFG	Asia	Afghanistan	04/01/20	NaN	0.0	NaN	NaN	0.0	NaN	...	NaN	59
...
57389	NaN	NaN	International	13/11/20	696.0	NaN	NaN	7.0	NaN	NaN	NaN	NaN	59
57390	NaN	NaN	International	14/11/20	696.0	NaN	NaN	7.0	NaN	NaN	NaN	NaN	59
57391	NaN	NaN	International	15/11/20	696.0	NaN	NaN	7.0	NaN	NaN	NaN	NaN	59
57392	NaN	NaN	International	16/11/20	696.0	NaN	NaN	7.0	NaN	NaN	NaN	NaN	59
57393	NaN	NaN	International	17/11/20	696.0	NaN	NaN	7.0	NaN	NaN	NaN	NaN	59

57394 rows × 50 columns

7. Data Aggregation:

a. Find max value in all columns using groupby function on 'continent' column Tip: use reset_index() after applying groupby.

```
In [8]: import pandas as pd
url = pd.read_csv("https://raw.githubusercontent.com/SR1608/Datasets/main/covid-data.csv")
grouped_df = url.groupby("continent")
maximums = grouped_df.max()
maximums = maximums.reset_index()
print(maximums)

continent iso_code location date \
0 Africa ZWE Zimbabwe 31/12/19
1 Asia YEM Yemen 31/12/19
2 Europe VAT Vatican 31/12/19
3 North America VIR United States Virgin Islands 31/12/19
4 Oceania WLF Wallis and Futuna 31/12/19
5 South America VEN Venezuela 31/12/19

0 total_cases new_cases new_cases_smoothed total_deaths new_deaths \
0 752269.0 13944.0 12583.714 20314.0 572.0
1 8874290.0 97894.0 93196.571 130519.0 2003.0
2 1301233.0 184813.0 156419.143 247220.0 4928.0
3 11205486.0 184813.0 156419.143 247220.0 4928.0
4 27750.0 1384.0 551.714 907.0 59.0
5 5876464.0 69074.0 46393.000 166014.0 3935.0

0 new_deaths_smoothed ... gdp_per_capita extreme_poverty \
0 297.429 ... 26382.287 77.6
1 1168.088 ... 16935.600 30.3
2 1101.008 ... 94277.965 5.7
3 2715.143 ... 54225.446 23.5
4 82.600 ... 44648.710 25.1
5 1096.714 ... 22767.037 7.1

0 cardiovascular_death diabetes_prevalence female_smokers male_smokers \
0 525.432 22.02 9.7 65.8
1 72.417 17.72 26.9 78.1
2 539.849 10.08 44.0 58.3
3 436.548 17.11 19.1 53.3
4 561.404 39.53 23.5 48.9
5 373.159 12.54 34.2 42.9

0 handwashing_facilities hospital_beds_per_thousand life_expectancy \
0 89.827 6.30 76.88
1 98.999 13.05 84.86
2 97.719 13.80 86.75
3 90.650 5.80 83.92
4 82.562 3.84 83.44
5 88.635 5.00 81.44

0 human_development_index
0 6.787
1 6.933
2 6.953
3 6.926
4 6.939
5 6.843

[6 rows x 49 columns]
```

b. Store the result in a new dataframe named 'df_groupby' (Use df_groupby dataframe for all further analysis)

```
In [12]: import pandas as pd
df_groupby=pd.DataFrame(maximums)
display(df_groupby)

Out[12]:
```

	continent	iso_code	location	date	total_cases	new_cases	new_cases_smoothed	total_deaths	new_deaths	new_deaths_smoothed	...	gdp_per_capita	extreme_poverty	cardiovascular_death
0	Africa	ZWE	Zimbabwe	31/12/19	752269.0	13944.0	12583.714	20314.0	572.0	297.429	...	77.6	59	525.432
1	Asia	YEM	Yemen	31/12/19</										