

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

1. pd.read_csv

pd.read_csv is a function in the pandas library in Python that is used to read a CSV (Comma Separated Values) file and convert it into a pandas DataFrame.

```
In [2]: df = pd.read_csv('C:/Users/ADMIN/Desktop/Python/Datasets/Covid Dataset/covid_worldwide.csv')
```

2. df.head()

This function will give you the first 5 rows of the dataframe. This will help you to get the simple idea about the \ dataset

```
In [3]: df.head()
```

	Serial Number	Country	Total Cases	Total Deaths	Total Recovered	Active Cases	Total Test	Population
0	1	USA	104196861	1132935	101322779	1741147	1159832679	334805269
1	2	India	44682784	530740	44150289	1755	915265788	1406631776
2	3	France	39524311	164233	39264546	95532	271490268	65584518
3	4	Germany	37798333	165711	37398100	216022	122332468	83883596
4	5	Brazil	36824580	697074	35919372	208134	63776166	215353593

3.df.isnull().sum()

This function will give you the info about which columns have the missing value in total

```
In [4]: df.isnull().sum()
```

	Serial Number	Country	Total Cases	Total Deaths	Total Recovered	Active Cases	Total Test	Population
Country	0							
Total Cases	0							
Total Deaths	0							
Total Recovered	21							
Active Cases	19							
Total Test	18							
Population	3							

dtype: int64

4.df.dropna

The dropna() method removes the rows that contains NULL values.

The dropna() method returns a new DataFrame object unless the inplace parameter is set to True, in that case the dropna() method does the removing in the original DataFrame instead.

```
In [5]: df.dropna()
```

	Serial Number	Country	Total Cases	Total Deaths	Total Recovered	Active Cases	Total Test	Population
0	1	USA	104196861	1132935	101322779	1741147	1159832679	334805269
1	2	India	44682784	530740	44150289	1755	915265788	1406631776
2	3	France	39524311	164233	39264546	95532	271490188	65584518
3	4	Germany	37798333	165711	37398100	216022	122332384	83883596
4	5	Brazil	36824580	697074	35919372	208134	63776166	215353593
...
217	218	Anguilla	3904	12	3879	13	51382	15230
218	219	Macao	3488	120	3357	11	7850	667490
219	220	Saint Pierre Miquelon	3452	2	2449	1001	25400	5759
220	221	Wallis and Futuna	3427	7	438	2982	20508	10982
224	225	Montserrat	1403	8	1376	19	17762	4965

195 rows × 8 columns

5.df.info()

This function gives you a basic idea about the dataset what are the type of the data, how many null value do we have \ in the data,how many no of features do we have , whats the memory usage.

```
In [6]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 231 entries, 0 to 230
Data columns (total 8 columns):
# Column      Non-Null Count  Dtype
---  ---
0 Serial Number      231 non-null    int64
1 Country            231 non-null    object
2 Total Cases        225 non-null    object
3 Total Deaths       225 non-null    object
4 Total Recovered    210 non-null    object
5 Active Cases       212 non-null    object
6 Total Test         213 non-null    object
7 Population         228 non-null    object
dtypes: int64(1), object(7)
memory usage: 14.6+ KB
```

6. df.astype

This function is very imp to change the datatype of the column. The reason for changing the data type is to do some numerical calculation

```
In [7]: df['Total Cases'] = df['Total Cases'].str.replace(',',' ', regex=True).astype("float")

In [8]: df['Total Deaths'] = df['Total Deaths'].str.replace(',',' ', regex=True).astype("float")

In [9]: df['Total Recovered'] = df['Total Recovered'].str.replace(',',' ', regex=True).astype("float")

In [10]: df['Active Cases'] = df['Active Cases'].str.replace(',',' ', regex=True).astype("float")

In [11]: df['Total Test'] = df['Total Test'].str.replace(',',' ', regex=True).astype("float")

In [12]: df['Population'] = df['Population'].str.replace(',',' ', regex=True).astype("float")

In [13]: df.info()
```

```
## as we can see we have change the data type from object to float now we will be able to perform some numerical problems
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 231 entries, 0 to 230
Data columns (total 8 columns):
# Column      Non-Null Count  Dtype
---  ---
0 Serial Number      231 non-null    int64
1 Country            231 non-null    object
2 Total Cases        231 non-null    float64
3 Total Deaths       225 non-null    float64
4 Total Recovered    210 non-null    float64
5 Active Cases       212 non-null    float64
6 Total Test         213 non-null    float64
7 Population         228 non-null    float64
dtypes: float64(6), int64(1), object(1)
memory usage: 14.6+ KB
```

7.df.describe()

The df.describe() method in pandas is used to generate summary statistics of various features of a DataFrame. It returns a new DataFrame that contains the count, mean, standard deviation, minimum, 25th percentile, median, 75th percentile, and maximum of each numerical column in the original DataFrame

```
In [14]: df.describe()
```

	Serial Number	Total Cases	Total Deaths	Total Recovered	Active Cases	Total Test	Population
count	231.000000	2.310000e+02	2.250000e+02	2.100000e+02	2.120000e+02	2.130000e+02	2.280000e+02
mean	116.000000	2.923400e+06	3.005778e+04	2.993905e+06	9.099855e+04	3.249315e+07	2.849325e+07
std	66.828138	9.478296e+06	1.053800e+05	9.520209e+06	7.663886e+05	1.173730e+08	1.022803e+08
min	1.000000	5.000000e+00	1.000000e+00	2.000000e+00	0.000000e+00	7.850000e+03	7.990000e+02
25%	58.500000	2.400100e+04	2.230000e+02	2.197250e+04	5.250000e+01	3.478150e+05	4.450815e+05
50%	116.000000	2.065920e+05	2.179000e+03	2.351455e+05	1.115000e+03	2.261693e+06	5.676382e+06
75%	173.500000	1.296146e+06	1.445200e+04	1.465768e+06	1.022875e+04	1.298426e+07	2.170759e+07
max	231.000000	1.041968e+08	1.132935e+06	1.013228e+08	1.095262e+07	1.159833e+09	1.406632e+09

You can also include or exclude certain columns and also include non-numerical columns by passing appropriate arguments to the method.

```
In [ ]: 
```

```
In [15]: df.describe(include = 'all')
```

	Serial Number	Country	Total Cases	Total Deaths	Total Recovered	Active Cases	Total Test	Population
count	231.000000	231	2.310000e+02	2.250000e+02	2.100000e+02	2.120000e+02	2.130000e+02	2.280000e+02
unique	NaN	231	NaN	NaN	NaN	NaN	NaN	NaN
top	NaN	USA	NaN	NaN	NaN	NaN	NaN	NaN
freq	NaN	1	NaN	NaN	NaN	NaN	NaN	NaN
mean	116.000000	NaN	2.923400e+06	3.005778e+04	2.993905e+06	9.099855e+04	3.249315e+07	2.849325e+07
std	66.828138	NaN	9.478296e+06	1.053800e+05	9.520209e+06	7.663886e+05	1.173730e+08	1.022803e+08
min	1.000000	NaN	5.000000e+00	1.000000e+00	2.000000e+00	0.000000e+00	7.850000e+03	7.990000e+02
25%	58.500000	NaN	2.400100e+04	2.230000e+02	2.197250e+04	5.250000e+01	3.478150e+05	4.450815e+05
50%	116.000000	NaN	2.065920e+05	2.179000e+03	2.351455e+05	1.115000e+03	2.261693e+06	5.676382e+06
75%	173.500000	NaN	1.296146e+06	1.445200e+04	1.465768e+06	1.022875e+04	1.298426e+07	2.170759e+07
max	231.000000	NaN	1.041968e+08	1.132935e+06	1.013228e+08	1.095262e+07	1.159833e+09	1.406632e+09

```
In [16]: df.describe(exclude = 'number') # to exclude the numerical column
```

	Country
count	231
unique	231
top	USA
freq	1

8.df.drop()

This function is usefull in dropping a particular column

```
In [17]: # serial nuber column is of no use so its better to drop the column
df.drop("Serial Number", axis = 1, inplace = True)
# here we have to sepcify the axis as i cause the function needs to know whether we are deleting a column or a row
# here inplace means the changes should be reflected in the original dataset

In [18]: df
```

```
# as we can see now there is no serial number column
```

	Country	Total Cases	Total Deaths	Total Recovered	Active Cases	Total Test	Population
0	USA	104196861.0	1132935.0	101322779.0	1741147.0	1.159833e+09	3.348053e+08
1	India	44682784.0	530740.0	44150289.0	1755.0	9.152658e+08	1.406632e+09
2	France	39524311.0	164233.0	39264546.0	95532.0	2.714902e+08	6.558452e+07
3	Germany	37798333.0	165711.0	37398100.0	216022.0	1.223324e+08	8.388350e+07
4	Brazil	36824580.0	697074.0	35919372.0	208134.0	6.377617e+07	2.153536e+08
...
226	Diamond Princess	712.0	13.0	699.0	0.0	NaN	NaN
227	Vatican City	29.0	NaN	29.0	0.0	NaN	7.990000e+02
228	Western Sahara	10.0	1.0	9.0	0.0	NaN	6.261510e+05
229	MS Zaandam	9.0	2.0	7.0	0.0	NaN	NaN
230	Tokelau	5.0	NaN	NaN	5.0	NaN	1.378000e+03

231 rows × 7 columns

9.df.isnull()

This function will help us to know which columns have NAN values

```
In [19]: df.isnull()
```

```
# this function will give a dataframe with true and false value, if the value is false then there is no null value
# id the value is true then there is a null value
```

	Country	Total Cases	Total Deaths	Total Recovered	Active Cases	Total Test	Population
0	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False
...
226	False	False	False	False	False	True	True
227	False	False	True	False	False	True	False
228	False	False	False	False	False	True	False
229	False	False	False	False	False	True	True
230	False	False	True	True	False	True	False

231 rows × 7 columns

10.The fillna()

method in Pandas is used to fill in missing values in a DataFrame with a specified value or method. By default, it replaces missing values with NaN, but you can specify a different value to use instead as shown below:

1.value: Specifies the value to use to fill in the missing values. Can be a scalar value or a dict of values for different columns. \ 2.method: Specifies the method to use for filling in missing values. Can be 'bfill' (backward-fill) or 'bfill' (backward-fill) or 'interpolate'(interpolate values) or 'pad' or 'backfill' \ 3.axis: Specifies the axis along which to fill in missing values. It can be 0 (rows) or 1 (columns). \ 4.inplace: Whether to fill in the missing values in place (modifying the original DataFrame) or to return a new DataFrame with the missing values filled in. \ 5.limit: Specifies the maximum number of consecutive missing values to fill. \ 6.downcast: Specifies a dictionary of values to use to downcast the data types of columns. \

```
In [20]: df['Total Deaths'].fillna(0, inplace = True)
# here we have filled the na value with 0

In [21]: df['Total Deaths'].isnull().sum()
# now as we can see we have 0 null value left same we can do for other column

In [22]: df
```

```
0
```

```
In [22]: df.dropna(inplace=True)

In [23]: df.isnull().sum()
# now by using the drop na function we have removed all the na values
```

	Country	Total Cases	Total Deaths	Total Recovered	Active Cases	Total Test	Population
Country	0						
Total Cases	0						
Total Deaths	0						
Total Recovered	0						
Active Cases	0						
Total Test	0						
Population	0						

dtype: int64

Now we have covered all the ways of dealing with missing value

10.df.assign()

pandas' assign() function is used to add new columns to a DataFrame, based on the computation of existing columns. It allows you to add new columns to a DataFrame without modifying the original dataframe. The function returns a new DataFrame with the added columns.

It's important to note that the original DataFrame df remains unchanged and the new DataFrame df_new is returned with the new columns added.

```
In [24]: df_new = df.assign(total_deaths_plus5 = df['Total Deaths']+5)
```

```
In [25]: df_new
# as you can see we have created a new column name total_deaths_plus_5
```

```
Out[25]:
```

	Country	Total Cases	Total Deaths	Total Recovered	Active Cases	Total Test	Population	total_deaths_plus5
0	USA	104196861.0	1132935.0	101322779.0	1741147.0	1.159833e+09	3.348053e+08	1132940.0
1	India	44682784.0	530740.0	44150289.0	1755.0	9.152658e+08	1.406632e+09	530745.0
2	France	39524311.0	164233.0	39264546.0	95532.0	2.714902e+08	6.558452e+07	164238.0
3	Germany	37798333.0	165711.0	37398100.0	216022.0	1.223324e+08	8.388350e+07	165716.0
4	Brazil	36824580.0	697074.0	35919372.0	208134.0	6.377617e+07	2.153536e+08	697079.0
...
218	Macao	3488.0	120.0	3357.0	11.0	7.850000e+03	6.674900e+05	125.0
219	Saint Pierre Miquelon	3452.0	2.0	2449.0	1001.0	2.540000e+04	5.759000e+03	7.0
220	Wallis and Futuna	3427.0	7.0	438.0	2982.0	2.050800e+04	1.098200e+04	12.0
223	Falkland Islands	1930.0	0.0	1930.0	0.0	8.632000e+03	3.539000e+03	5.0
224	Montserrat	1403.0	8.0	1376.0	19.0	1.776200e+04	4.965000e+03	13.0

196 rows × 8 columns

11.df.columns()

In the above method you have to provide a list of column names

```
In [26]: df.columns = ['Country', 'Total_cases', 'Total_deaths', 'Total_recovered', 'Active_cases', 'Total_test', 'Population']

In [27]: df
```

```
Out[27]:
```

	Country	Total_cases	Total_deaths	Total_recovered	Active_cases	Total_test	Population
0	USA	104196861.0	1132935.0	101322779.0	1741147.0	1.159833e+09	3.348053e+08
1	India	44682784.0	530740.0	44150289.0	1755.0	9.152658e+08	1.406632e+09
2	France	39524311.0	164233.0	39264546.0	95532.0	2.714902e+08	6.558452e+07
3	Germany	37798333.0	165711.0	37398100.0	216022.0	1.223324e+08	8.388350e+07
4	Brazil	36824580.0	697074.0	35919372.0	208134.0	6.377617e+07	2.153536e+08
...
218	Macao	3488.0	120.0	3357.0	11.0	7.850000e+03	6.674900e+05
219	Saint Pierre Miquelon	3452.0	2.0	2449.0	1001.0	2.540000e+04	5.759000e+03
220	Wallis and Futuna	3427.0	7.0	438.0	2982.0	2.050800e+04	1.098200e+04
223	Falkland Islands	1930.0	0.0	1930.0	0.0	8.632000e+03	3.539000e+03
224	Montserrat	1403.0	8.0	1376.0	19.0	1.776200e+04	4.965000e+03

196 rows × 7 columns

What if you want to assign a specific name to particular column

12.df.query()

Pandas' query() function allows you to filter a DataFrame based on a Boolean expression. It allows you to select rows from a DataFrame using a query string similar to SQL.\

The function returns a new DataFrame containing only the rows that satisfy the Boolean expression

```
In [28]: df_query = df.query('Total_deaths > 1000')
df_query
```

|--|