

# install and import packages to the enviroment

In [1]:

```

1  #install packages required for the enviroment.
2  #!pip install chart_studio
3  #!pip install autoviz
4  #!pip install xlrd
5  #!pip install -q sklearn
6
7  #import packages to the enviroment
8  import pandas as pd
9  import numpy as np
10 import chart_studio.plotly as py
11 import cufflinks as cf
12 import seaborn as sns
13 import plotly.express as px
14 import plotly.graph_objects as go
15 from autoviz.AutoViz_Class import AutoViz_Class
16 import h2o
17 from h2o.automl import H2OAutoML
18 import pandas as pd
19 h2o.init()
20
21
22 import tensorflow as tf
23
24 from tensorflow import feature_column
25 from tensorflow.keras import layers
26 from sklearn.model_selection import train_test_split
27
28
29 %matplotlib inline
30
31 # Make Plotly work in your Jupyter Notebook
32 from plotly.offline import download_plotlyjs, init_notebook_mode, plot, iplot
33 init_notebook_mode(connected=True)
34 # Use Plotly Locally
35 cf.go_offline()
36
37 print("Packages installed in enviroment successfully")

```

Imported AutoViz\_Class version: 0.0.81. Call using:

```

from autoviz.AutoViz_Class import AutoViz_Class
AV = AutoViz_Class()
AV.AutoViz(filename, sep=',', depVar='', dfte=None, header=0, verbose=
0,

```

```

        lowess=False, chart_format='svg', max_rows_analy
zed=150000, max_cols_analyzed=30)

```

Note: verbose=0 or 1 generates charts and displays them in your local Jupyter notebook.

verbose=2 saves plots in your local machine under AutoViz\_Plots directory and does not display charts.

Checking whether there is an H2O instance running at <http://localhost:54321> (<http://localhost:54321>) ..... not found.

Attempting to start a local H2O server...

; Java HotSpot(TM) 64-Bit Server VM (build 25.221-b11, mixed mode)

Starting server from C:\Users\MSI\anaconda3\Lib\site-packages\h2o\backend\bin\h2o.jar

Ice root: C:\Users\MSI\AppData\Local\Temp\tmph8f1847m

JVM stdout: C:\Users\MSI\AppData\Local\Temp\tmph8f1847m\h2o\_MSI\_started\_from\_python.out

JVM stderr: C:\Users\MSI\AppData\Local\Temp\tmph8f1847m\h2o\_MSI\_started\_

```
from_python.err
```

```
Server is running at http://127.0.0.1:54321 (http://127.0.0.1:54321)
Connecting to H2O server at http://127.0.0.1:54321 (http://127.0.0.1:54321) ... successful.
```

```
H2O_cluster_uptime:          02 secs
H2O_cluster_timezone:        Europe/Paris
H2O_data_parsing_timezone:    UTC
H2O_cluster_version:          3.32.1.3
H2O_cluster_version_age:      3 days
H2O_cluster_name: H2O_from_python_MSI_e7bj7w
H2O_cluster_total_nodes:      1
H2O_cluster_free_memory:      1.766 Gb
H2O_cluster_total_cores:      4
H2O_cluster_allowed_cores:    4
H2O_cluster_status: accepting new members,
                    healthy
H2O_connection_url: http://127.0.0.1:54321
H2O_connection_proxy: {"http": null, "https": null}
H2O_internal_security: False
H2O_API_Extensions: Amazon S3, Algos,
                    AutoML, Core V3,
                    TargetEncoder, Core V4
Python_version: 3.8.5 final
```

Packages installed in enviroment successfully

## Load Dataset and inspect

Now that we have installed and imported the packages that we will use for the enviroment. lets go ahead and load the dataset for inspection and analysis.

In [2]:

```

1 #Load the dataset
2 df = pd.read_csv("country_wise_latest.csv")
3
4 print("'Dataset information'")
5 print()
6 df.info()

```

'Dataset information'

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 187 entries, 0 to 186
Data columns (total 15 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Country/Region                        187 non-null    object
1   Confirmed                            187 non-null    int64
2   Deaths                              187 non-null    int64
3   Recovered                           187 non-null    int64
4   Active                              187 non-null    int64
5   New cases                           187 non-null    int64
6   New deaths                          187 non-null    int64
7   New recovered                        187 non-null    int64
8   Deaths / 100 Cases                  187 non-null    float64
9   Recovered / 100 Cases                187 non-null    float64
10  Deaths / 100 Recovered              187 non-null    float64
11  Confirmed last week                  187 non-null    int64
12  1 week change                        187 non-null    int64
13  1 week % increase                    187 non-null    float64
14  WHO Region                          187 non-null    object
dtypes: float64(4), int64(9), object(2)
memory usage: 22.0+ KB

```

In [3]:

```

1 #columns in the dataset
2 df.columns

```

Out[3]:

```

Index(['Country/Region', 'Confirmed', 'Deaths', 'Recovered', 'Active',
      'New cases', 'New deaths', 'New recovered', 'Deaths / 100 Cases',
      'Recovered / 100 Cases', 'Deaths / 100 Recovered',
      'Confirmed last week', '1 week change', '1 week % increase',
      'WHO Region'],
      dtype='object')

```

In [4]:

```
1 #describe the dataset
2 df.describe()
```

Out[4]:

	Confirmed	Deaths	Recovered	Active	New cases	New deaths
count	1.870000e+02	187.000000	1.870000e+02	1.870000e+02	187.000000	187.000000
mean	8.813094e+04	3497.518717	5.063148e+04	3.400194e+04	1222.957219	28.957219
std	3.833187e+05	14100.002482	1.901882e+05	2.133262e+05	5710.374790	120.037173
min	1.000000e+01	0.000000	0.000000e+00	0.000000e+00	0.000000	0.000000
25%	1.114000e+03	18.500000	6.265000e+02	1.415000e+02	4.000000	0.000000
50%	5.059000e+03	108.000000	2.815000e+03	1.600000e+03	49.000000	1.000000
75%	4.046050e+04	734.000000	2.260600e+04	9.149000e+03	419.500000	6.000000
max	4.290259e+06	148011.000000	1.846641e+06	2.816444e+06	56336.000000	1076.000000



Lets examine the top 5 countries with max confirmed numbers.

In [5]:

```

1  #max number of confirmed, recovered, deaths and active.
2  maxConfirmed = df.loc[df["Confirmed"].idxmax()]
3  maxRecovered = df.loc[df["Recovered"].idxmax()]
4  maxDeaths = df.loc[df["Deaths"].idxmax()]
5  maxActive = df.loc[df["Active"].idxmax()]
6
7  print("-----Max Confirmed number-----")
8  print(maxConfirmed)
9  print()
10 print("-----Max Deaths number-----")
11 print(maxDeaths)
12 print()
13 print("-----Max Recovered number-----")
14 print(maxRecovered)
15 print()
16 print("-----Max Active number-----")
17 print(maxActive)

```

-----Max Confirmed number-----

Country/Region	US
Confirmed	4290259
Deaths	148011
Recovered	1325804
Active	2816444
New cases	56336
New deaths	1076
New recovered	27941
Deaths / 100 Cases	3.45
Recovered / 100 Cases	30.9
Deaths / 100 Recovered	11.16
Confirmed last week	3834677
1 week change	455582
1 week % increase	11.88
WHO Region	Americas

Name: 173, dtype: object

-----Max Deaths number-----

Country/Region	US
Confirmed	4290259
Deaths	148011
Recovered	1325804
Active	2816444
New cases	56336
New deaths	1076
New recovered	27941
Deaths / 100 Cases	3.45
Recovered / 100 Cases	30.9
Deaths / 100 Recovered	11.16
Confirmed last week	3834677
1 week change	455582
1 week % increase	11.88
WHO Region	Americas

Name: 173, dtype: object

-----Max Recovered number-----

Country/Region	Brazil
Confirmed	2442375
Deaths	87618
Recovered	1846641

```

Active          508116
New cases       23284
New deaths      614
New recovered   33728
Deaths / 100 Cases  3.59
Recovered / 100 Cases  75.61
Deaths / 100 Recovered  4.74
Confirmed last week  2118646
1 week change    323729
1 week % increase  15.28
WHO Region       Americas
Name: 23, dtype: object

```

```

-----Max Active number-----
Country/Region      US
Confirmed           4290259
Deaths              148011
Recovered           1325804
Active              2816444
New cases            56336
New deaths           1076
New recovered        27941
Deaths / 100 Cases   3.45
Recovered / 100 Cases  30.9
Deaths / 100 Recovered  11.16
Confirmed last week  3834677
1 week change        455582
1 week % increase    11.88
WHO Region           Americas
Name: 173, dtype: object

```

Lets visualise the top 10 confirmed countries in descending order

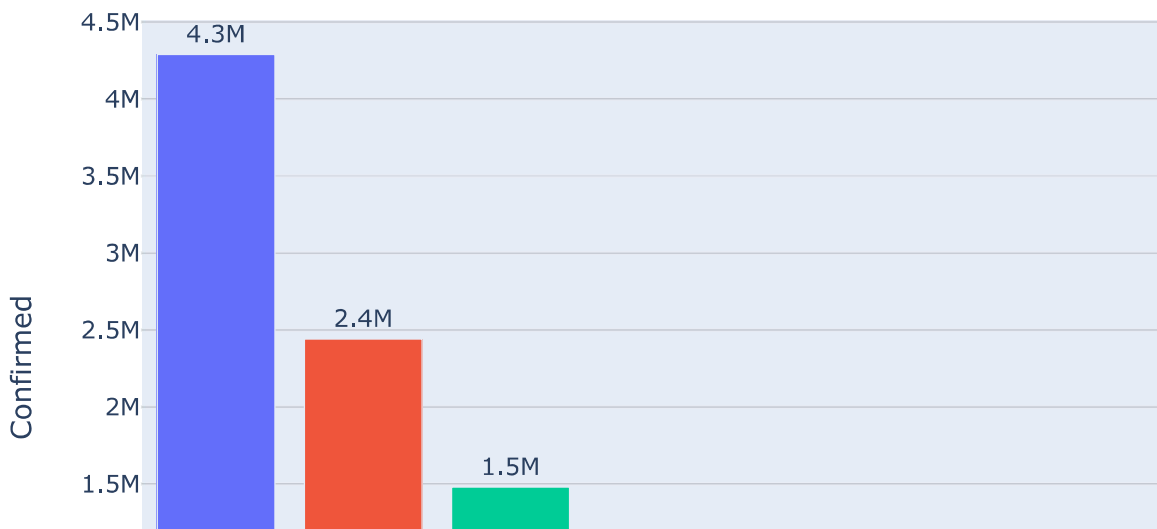
In [6]:

```

1 #dislay confirmed in first 10 countries
2 df = df.sort_values("Confirmed", ascending = False)
3 fig = px.bar(df.head(10), y = "Confirmed", x = "Country/Region",
4             text = "Confirmed", color = "Country/Region")
5 #total values bar with 2 precision values
6 fig.update_traces(texttemplate = "%{text:.2s}", textposition = "outside")
7 #set fontsize and uniformText
8 fig.update_layout(uniformtext_minsize = 8)
9 #rotate Label 45 degrees
10 fig.update_layout(xaxis_tickangle=-45)
11 fig.update_layout(legend_title_text='confirmed (Covid-19)')
12 fig.update_layout(title_text='Top 10 Confirmed cases in the world')
13 fig
14

```

Top 10 Confirmed cases in the world



Top 10 Deaths reported countries in the world



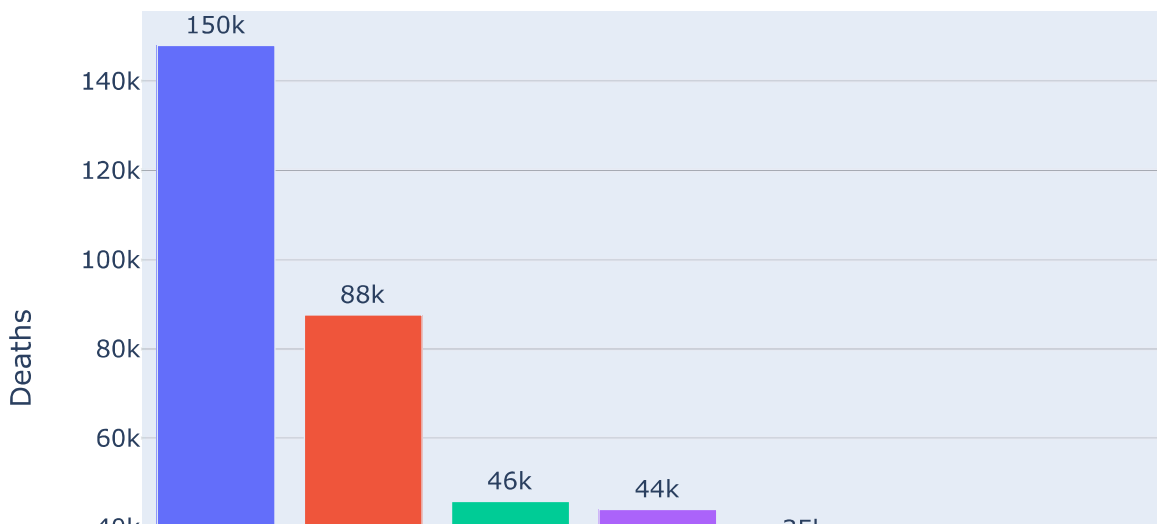
In [7]:

```

1 #display deaths in first 10 countries
2 df = df.sort_values("Deaths", ascending = False)
3 fig = px.bar(df.head(10), y = "Deaths", x = "Country/Region",
4             text = "Deaths", color = "Country/Region")
5 #total values bar with 2 precision values
6 fig.update_traces(texttemplate = "%{text:.2s}", textposition = "outside")
7 #set fontsize and uniformText
8 fig.update_layout(uniformtext_minsize = 8)
9 #rotate Label 45 degrees
10 fig.update_layout(xaxis_tickangle=-45)
11 fig.update_layout(legend_title_text='Deaths (Covid-19)')
12 fig.update_layout(title_text='Top 10 Death cases in the world')
13 fig

```

Top 10 Death cases in the world



Top 10 recovered countries in the world

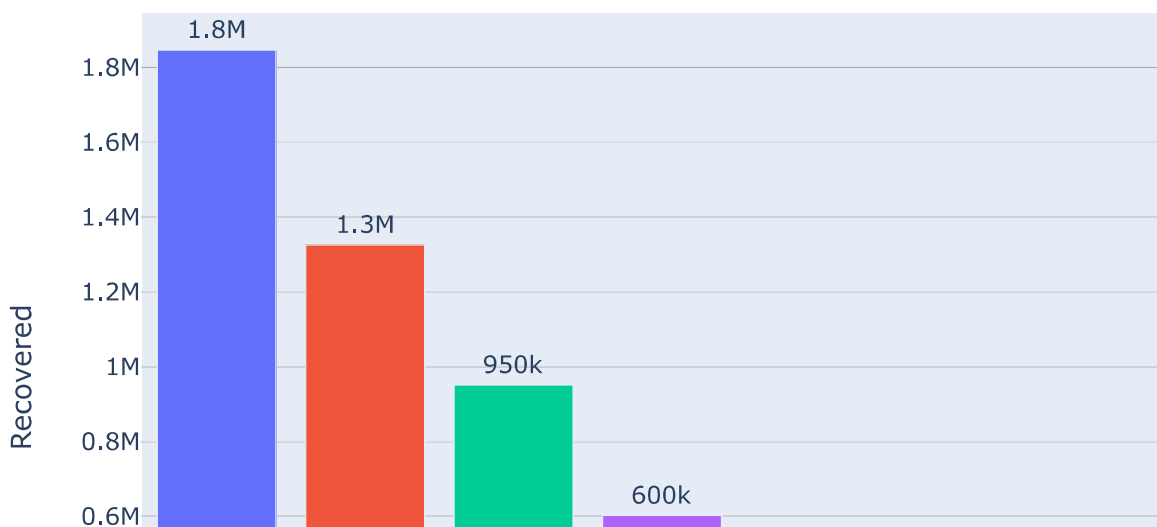
In [8]:

```

1 #dislay Recovered in first 10 countries
2 df = df.sort_values("Recovered", ascending = False)
3 fig = px.bar(df.head(10), y = "Recovered", x = "Country/Region",
4             text = "Recovered", color = "Country/Region")
5 #total values bar with 2 precision values
6 fig.update_traces(texttemplate = "%{text:.2s}", textposition = "outside")
7 #set fontsize and uniformText
8 fig.update_layout(uniformtext_minsize = 8)
9 #rotate Label 45 degrees
10 fig.update_layout(xaxis_tickangle=-45)
11 fig.update_layout(legend_title_text='Recovered (Covid-19)')
12 fig.update_layout(title_text='Top 10 Recovered cases in the world')
13 fig

```

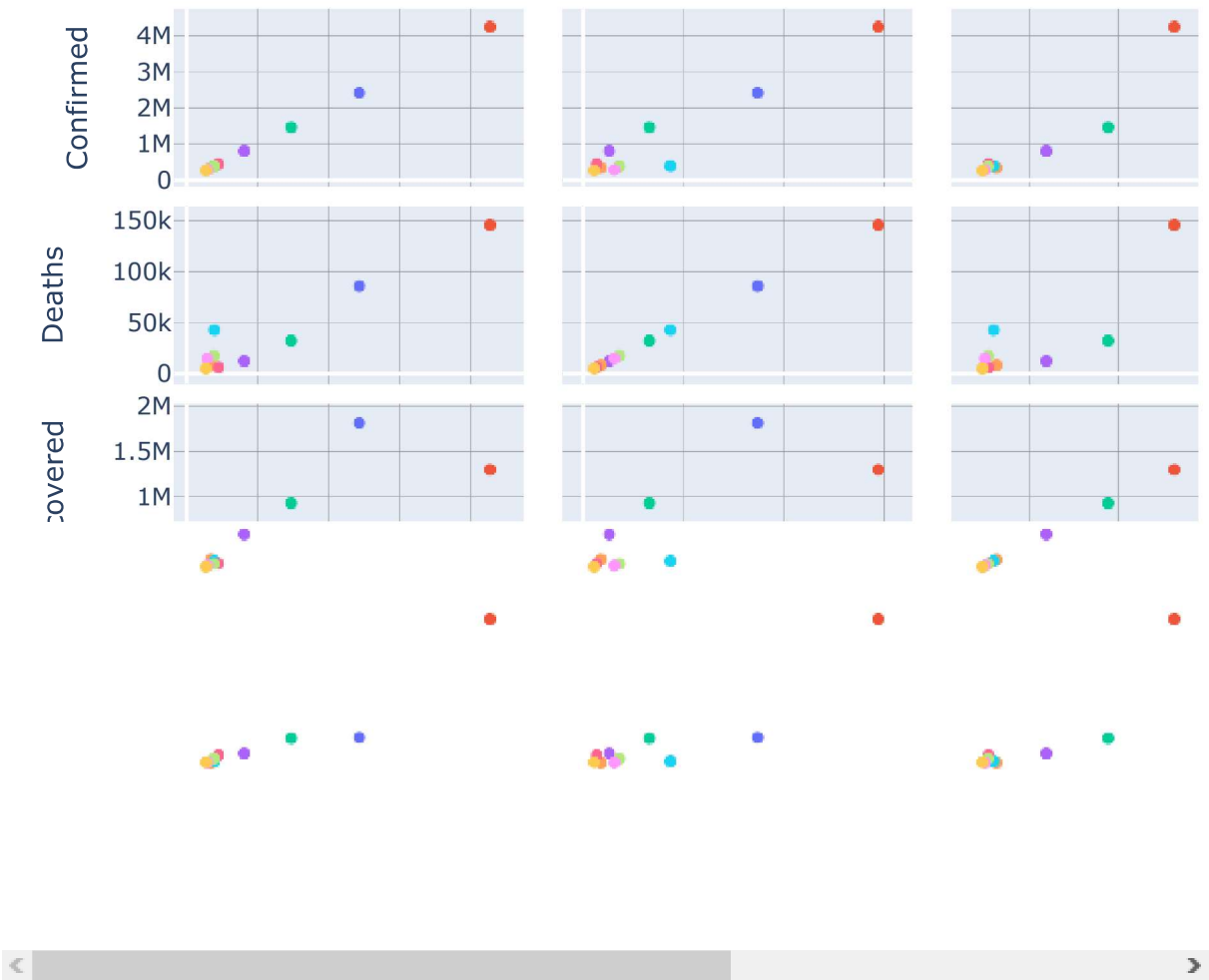
Top 10 Recovered cases in the world



Scatter matrix plot between the confirmed, deaths, recovered and active cases

In [9]:

```
1 #Scatter matrix plot between the confirmed, deaths and recovery category.
2 fig = px.scatter_matrix(df.head(10), dimensions=["Confirmed","Deaths","Recovered","Active"])
3 fig
```



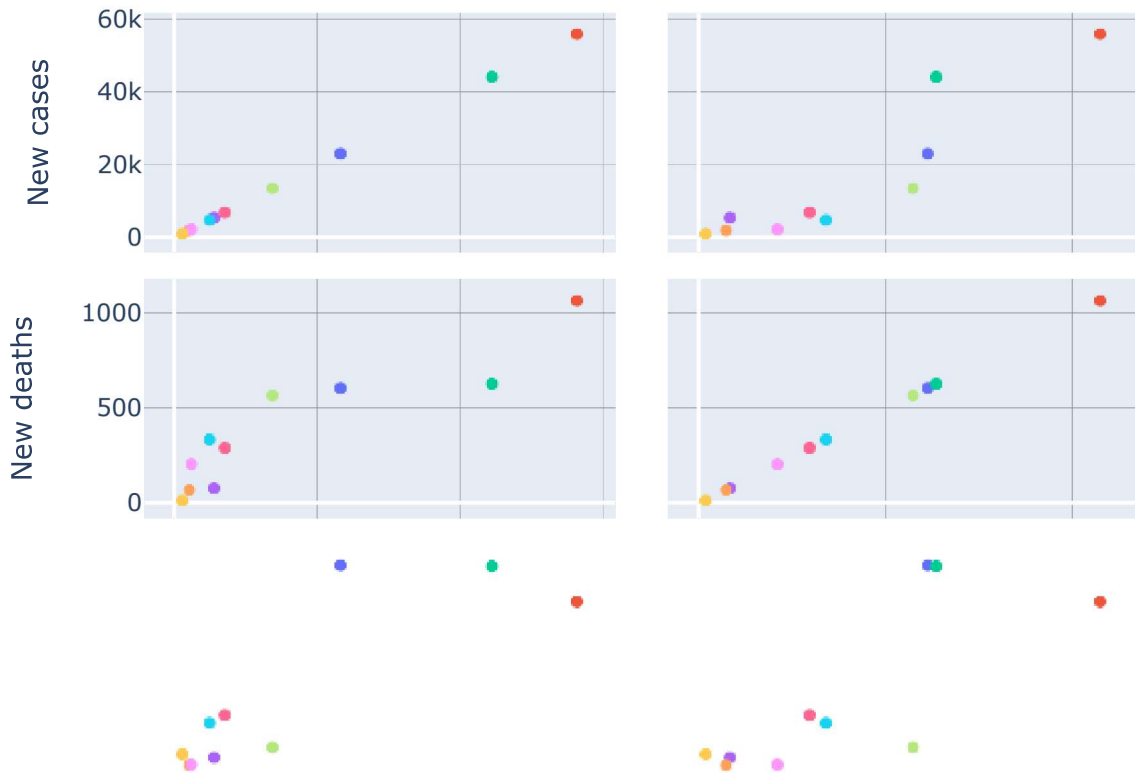
Scatter matrix of new cases, new deaths and new recovery.

In [10]:

```

1 #scatter matrix of new cases, new deaths and new recovery.
2 fig = px.scatter_matrix(df.head(10), dimensions = ["New cases", "New deaths", "New recovery"])
3 fig

```



Automatic visualisation of the entire data to find hidden patterns or insight.

In [11]:

```

1 av = AutoViz_Class()
2 df = "country_wise_latest.csv"
3 sep = ","
4 draftAutoViz = av.AutoViz(
5     df,
6     sep = ",",
7     depVar="",
8     dfte = None,
9     header = 0,
10    verbose = 0,
11    lowess = False,
12    chart_format = "svg",
13    max_cols_analyzed=30,
14    max_rows_analyzed=1500000,
15 )

```

Shape of your Data Set: (187, 15)

##### C L A S S I F Y I N G V A R I A B L E S #####  
 ###

Classifying variables in data set...

```

Number of Numeric Columns = 4
Number of Integer-Categorical Columns = 9
Number of String-Categorical Columns = 1
Number of Factor-Categorical Columns = 0
Number of String-Boolean Columns = 0
Number of Numeric-Boolean Columns = 0
Number of Discrete String Columns = 0
Number of NLP String Columns = 0
Number of Date Time Columns = 0
Number of ID Columns = 1
Number of Columns to Delete = 0
15 Predictors classified...

```

This does not include the Target column(s)

1 variables removed since they were ID or low-information variable

S

## create a prediction

In [12]:

```

1 df = pd.read_csv("country_wise_latest.csv")
2
3 df = df.sort_values("Confirmed", ascending = False)
4 df

```

Out[12]:

	Country/Region	Confirmed	Deaths	Recovered	Active	New cases	New deaths	New recovered	Death / 100 Case
173	US	4290259	148011	1325804	2816444	56336	1076	27941	3.4
23	Brazil	2442375	87618	1846641	508116	23284	614	33728	3.5
79	India	1480073	33408	951166	495499	44457	637	33598	2.2
138	Russia	816680	13334	602249	201097	5607	85	3077	1.6
154	South Africa	452529	7067	274925	170537	7096	298	9848	1.5
...	...	...	...	...	...	...	...	...	...
49	Dominica	18	0	18	0	0	0	0	0.0
140	Saint Kitts and Nevis	17	0	15	2	0	0	0	0.0
68	Greenland	14	0	13	1	1	0	0	0.0
75	Holy See	12	0	12	0	0	0	0	0.0
183	Western Sahara	10	1	8	1	0	0	0	10.0

187 rows × 15 columns

In [13]:

```

1 #split dataset to train and test dataset
2 x_train, y_train = train_test_split(df, test_size=0.2)
3 x_test, y_test = train_test_split(df, test_size = 0.2)
4 print(len(x_train), "x_train Dataset")
5 print(len(y_train), "y_train Dataset")
6 print(len(x_test), "x_test Dataset")
7 print(len(y_test), "y_test Dataset")

```

```

149 x_train Dataset
38 y_train Dataset
149 x_test Dataset
38 y_test Dataset

```

In [14]:

```

1 x_train = np.array(x_train)
2 y_train = np.array(y_train)
3 x_test = np.array(x_test)
4 y_test = np.array(y_test)
5
6 x_train

```

Out[14]:

```

array([[ 'Guinea-Bissau', 1954, 26, ..., 5, 0.26, 'Africa'],
       [ 'Uzbekistan', 21209, 121, ..., 4060, 23.67, 'Europe'],
       [ 'Namibia', 1843, 8, ..., 499, 37.13, 'Africa'],
       ...,
       [ 'Pakistan', 274289, 5842, ..., 8193, 3.08,
         'Eastern Mediterranean'],
       [ 'Poland', 43402, 1676, ..., 3019, 7.48, 'Europe'],
       [ 'Czechia', 15516, 373, ..., 1418, 10.06, 'Europe']], dtype=object)

```

In [15]:

```

1 df = h2o.import_file("country_wise_latest.csv")
2
3 df2 = df.sort(1, ascending=False)
4 df2

```

Parse progress:  100%

Country/Region	Confirmed	Deaths	Recovered	Active	New cases	New deaths	New recovered	Dea / Ca
US	4.29026e+06	148011	1.3258e+06	2.81644e+06	56336	1076	27941	3
Brazil	2.44238e+06	87618	1.84664e+06	508116	23284	614	33728	3
India	1.48007e+06	33408	951166	495499	44457	637	33598	2
Russia	816680	13334	602249	201097	5607	85	3077	1
South Africa	452529	7067	274925	170537	7096	298	9848	1
Mexico	395489	44022	303810	47657	4973	342	8588	1
Peru	389717	18418	272547	98752	13756	575	4697	4
Chile	347923	9187	319954	18782	2133	75	1859	2
United Kingdom	301708	45844	1437	254427	688	7	3	15
Iran	293606	15912	255144	22550	2434	212	1931	5

Out[15]:

In [16]:

```
1 train,test,valid = df2.split_frame(ratios=[.7, .15])
2
3 x = train.columns
4 y = "Confirmed"
5 x.remove(y)
6
7
8
9 train[y] = train[y].asfactor()
10 test[y] = test[y].asfactor()
11
12 aml = H2OAutoML(max_models=20, seed = 1)
13 aml.train(x=x,y=y,training_frame=train)
```

AutoML progress: |

15:11:24.398: AutoML: XGBoost is not available; skipping it.

████████████████████

15:13:07.440: Skipping training of model GBM\_5\_AutoML\_20210523\_151124 due to exception: water.exceptions.H2OModelBuilderIllegalArgumentException: Illegal argument(s) for GBM model: GBM\_5\_AutoML\_20210523\_151124. Details: ERRR on field: \_min\_rows: The dataset size is too small to split for min\_rows=100.0: must have at least 200.0 (weighted) rows, but have only 119.0.

██

15:19:34.297: StackedEnsemble\_BestOfFamily\_AutoML\_20210523\_151124 [StackedEnsemble best (built using top model from each algorithm type)] failed: java.lang.RuntimeException: java.lang.ArrayIndexOutOfBoundsException: 591

15:19:36.331: StackedEnsemble\_AllModels\_AutoML\_20210523\_151124 [StackedEnsemble all (built using all AutoML models)] failed: java.lang.RuntimeException: java.lang.ArrayIndexOutOfBoundsException: 2361

██████████ | 100%



```
1 lb = aml.leaderboard
2 lb.head(rows=lb.nrows)
```

```
In [18]:
1 test = test.drop("Confirmed")
2 test
3 predictions = aml.leader.predict(test)
```

In [19]:

```
1 predictions.describe()
```

Rows:36

Cols:119

predict		p10	p17	p18
type	enum	real	real	real
mins		0.00016353909202609677	0.00661551470655862	0.0002836714388956398
mean		0.008649500961384735	0.008102108121068645	0.009856949869950662
maxs		0.010652713094576738	0.013340946614035601	0.014392778655689433
sigma		0.0018200699434186916	0.00115351036528011	0.0023084134086096297
zeros		0	0	0
missing	0	0	0	0
0	45309	0.00016353909202609677	0.013340946614035601	0.0002836714388956398
1	18	0.005748725093107088	0.00762718067972786	0.014392778655689433
2	4548	0.007320875055427739	0.007756135180213147	0.013214322729023403
3	701	0.010257462572932993	0.0073860848428042815	0.009371200650917172
4	4599	0.009719109409275868	0.008277350141375952	0.009833570662447414
5	10621	0.009993711739429931	0.007608825868532975	0.012807158824680637
6	10621	0.008244337801454116	0.007372844595787201	0.012212652273807578
7	4599	0.010652713094576738	0.007855823839581484	0.010510809189887647
8	4548	0.0092108850191359	0.008089744863513727	0.010720320590857613
9	701	0.009904732945779982	0.007407866049795423	0.009093155854929414

