

Spread Visualization and Prediction of the COVID-19 Using Machine Learning and Deep learning.

Mohammad Faisal Danish

Project Proposal:-

The current destructive pandemic of coronavirus disease 2019 (COVID-19), caused by severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), was first reported in Wuhan, China, in December 2019. The outbreak has affected millions of people around the world and the number of infections and mortalities has been growing at an alarming rate. In such a situation, forecasting and proper study of the pattern of disease spread can inspire design better strategies to make more efficient decisions. Moreover, such studies play an important role in achieving accurate predictions.

Machine learning has numerous tools that can be used for visualization and prediction, and nowadays it is used worldwide for study of the pattern of COVID-19 spread, e.g. One of the main focus of the study in this project is to use machine learning techniques to analyze and visualize the spreading of the virus country-wise as well as globally during a specific period of time by considering confirmed cases, recovered cases and fatalities. The global impact of the outbreak on various aspects of life has been the focus of many studies, e.g. On the other hand, a pandemic can be forecast by considering a variety of parameters such as the impact of environmental factors, quarantine, age, gender and a lot more. The forecasting accuracy depends on the availability of proper data to base its predictions and provide an estimate of uncertainty. A challenge to use machine learning techniques for the current outbreak is that the datasets are not yet standardized by any standardization organization and the statistical anomalies are not considered. Also, the appropriate selection of parameters and the selection of the best machine learning model for prediction are other challenges involved in training a model. In this project, we are going to perform linear regression, Support vector machine, Multilayer perceptron, Ensemble methods, Time series with ARIMA and Prophet etc., on the Johns Hopkins University's COVID-19 data to anticipate the future effects of COVID-19 pandemic in India and some other countries. Moreover, we are going to study the impact of some parameters such as geographic conditions, economic statistics, population statistics, life expectancy, etc. in prediction of COVID-19 spread.

Introduction:

The current destructive pandemic of coronavirus disease 2019 (COVID-19), caused by severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) [13], was first reported in Wuhan, China, in December 2019. The outbreak has affected millions of people around the world and the number of infections and mortalities has been growing at an alarming rate. As of date, confirmed COVID-19 cases are more than 20.6 Cr in almost all countries. In such a situation, forecasting and proper study of the pattern of disease spread can inspire design better strategies to make more efficient decisions. Moreover, such studies play an important role in achieving accurate predictions.

Machine learning has numerous tools that can be used for visualization and prediction, and nowadays it is used worldwide to study the pattern of COVID-19 spread. One of the main focus of the study in this project is to use machine learning techniques to analyze and visualize the spreading of the virus country-wise as well as globally during a specific period of time by considering confirmed cases, recovered cases and fatalities. The global impact of the outbreak on various aspects of life has been the focus of many studies. On the other hand, a pandemic can be forecast by considering a variety of parameters such as the impact of environmental factors, quarantine, age, gender and a lot more.

The forecasting accuracy depends on the availability of proper data to base its predictions and provide an estimate of uncertainty. A challenge to use machine learning techniques for the current outbreak is that the datasets are not yet standardized by any standardization organization and the statistical anomalies are not considered. Also, the appropriate selection of parameters and the selection of the best machine learning model for prediction are other challenges involved in training a model. In this project, we are going to perform Linear regression, Support vector machine, Ensemble methods, Multilayer perceptron, Recurrent neural network-LSTM, ARIMA and Prophet, etc., on the Johns Hopkins University's COVID-19 data to anticipate the future effects of COVID-19 pandemic in the world, India and some other countries. Moreover, we are going to study the impact of some other parameters such as environmental factors, life expectancy, population statistics, etc., in prediction of COVID-19 spread.

2 Experimental data and results:-

The data is provided by the Johns Hopkins University Center for Systems Science and Engineering and contains three time series with the number of reported daily confirmed cases, recovered cases and deaths by country. This dataset is updated automatically on daily basis. In this project we employed data from 22 January 2020 up to 7 July 2021. Initially, data preprocessing was almost challenging and much time was required because the dataset was not standard and many data cleaning processes were required. This part was done carefully and some appropriate data frames were prepared, such as follows.

	Date	Country/Region	Confirmed	Deaths	Recovered	Active	New confirmed	New deaths	New recovered	WHO region
34403	2020-07-23	West Bank and Gaza	9744	67	2720	6957	346	1	770	EMRO
34404	2020-07-23	Western Sahara	10	1	8	1	0	0	0	AFRO
34405	2020-07-23	Yemen	1654	461	762	431	14	3	11	EMRO
34406	2020-07-23	Zambia	3789	134	1677	1978	206	6	0	AFRO
34407	2020-07-23	Zimbabwe	2124	28	510	1586	90	2	0	AFRO

	Country/Region	Confirmed	Deaths	Recovered	Active	New confirmed	New deaths	New recovered	Recovery rate(per 100)	Mortality rate(per 100)	WHO region
0	Afghanistan	35928	1211	24550	10167	201	21	626	68.33	3.37	EMRO
1	Albania	4466	123	2523	1820	108	3	60	56.49	2.75	EURO
2	Algeria	25484	1124	17369	6991	612	13	386	68.16	4.41	AFRO
3	Andorra	889	52	803	34	0	0	0	90.33	5.85	EURO
4	Angola	851	33	236	582	39	0	15	27.73	3.88	AFRO

	Date	Confirmed	Deaths	Recovered	Active	New confirmed	New deaths	New recovered	Recovery rate(per 100)	Mortality rate(per 100)	Number of countries
0	2020-01-22	555	17	28	510	0	0	0	5.05	3.06	6
1	2020-01-23	654	18	30	606	99	1	2	4.59	2.75	8
2	2020-01-24	941	26	36	879	287	8	6	3.83	2.76	9
3	2020-01-25	1434	42	39	1353	493	16	3	2.72	2.93	11
4	2020-01-26	2118	56	52	2010	684	14	13	2.46	2.64	13

After exploring the data, we performed some visualizations on the data in order to get a better understanding of the data and how the pandemic is affecting all of us. For example, in Figure 1, we can see the latest status of cases in the world.

◆	Confirmed ◆	Recovered ◆	Deaths ◆	Active ◆	Recovery rate(per 100) ◆	Mortality rate(per 100) ◆	Number of countries ◆
554	196625810.00	127565256.00	4198924.00	64861630.00	-0.65	2.14	193.00

Figure 1 :- Total status of the COVID-19 Cases of the world

Also in Figure 2, we can see that the latest global recovery rate per 100 cases is 64.88% whereas the mortality rate per 100 cases is 2.14 that is a good news because at the start point of this project, the recovery rate was around 54 percent whereas the mortality rate was around 3 percent.

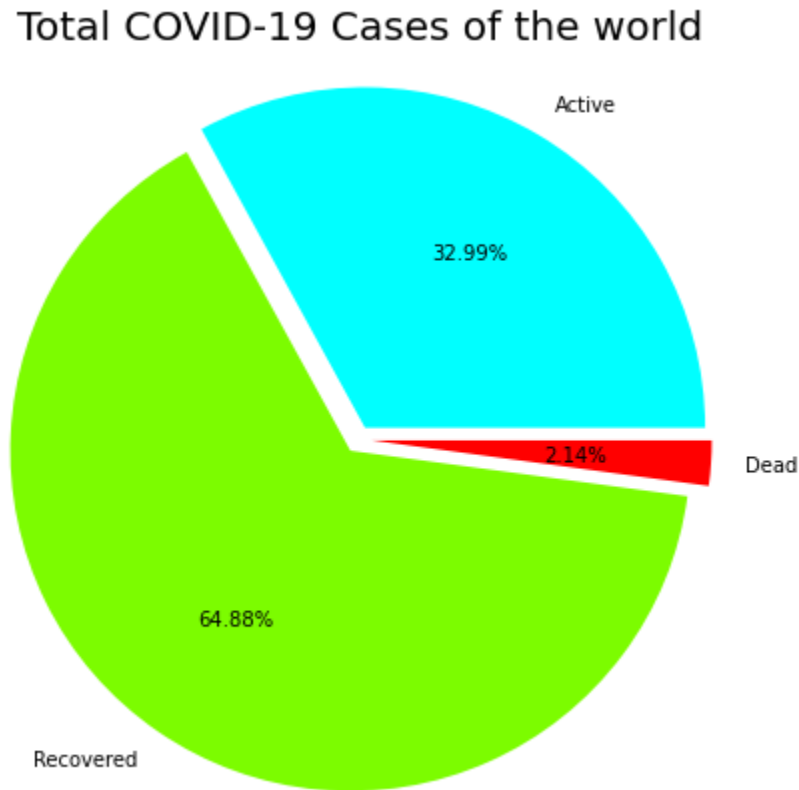
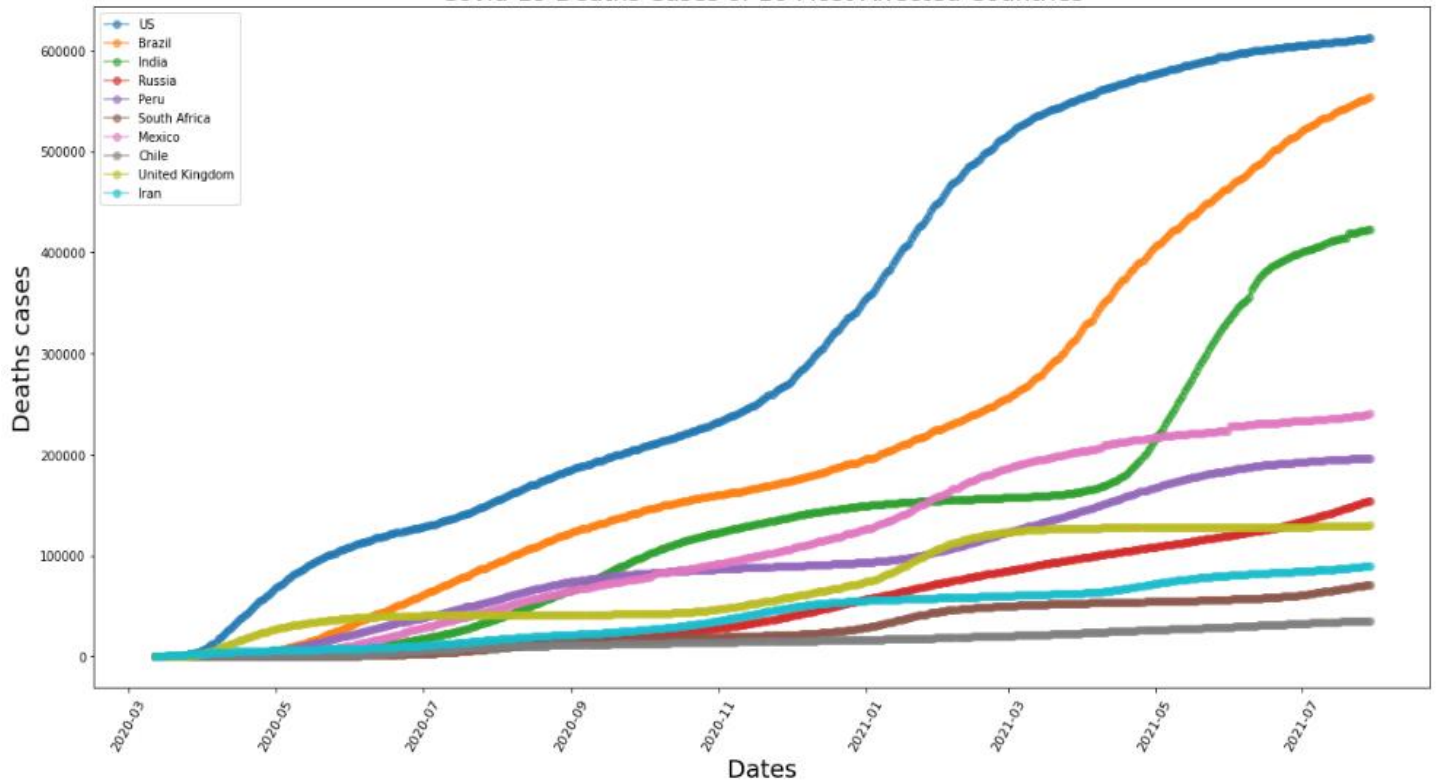


Figure 2:- Total Cases World Wise

Also as an example, Figure 3 shows comparisons between the latest COVID-19 cases status of 10 most affected countries, i.e., US, Brazil, India, Russia, Peru, South Africa, Mexico, Chile, United Kingdom, and Iran. Some of our conclusions based on the analysis from the above observations and some others, which can be found in the project's GitHub repository, are as follows:

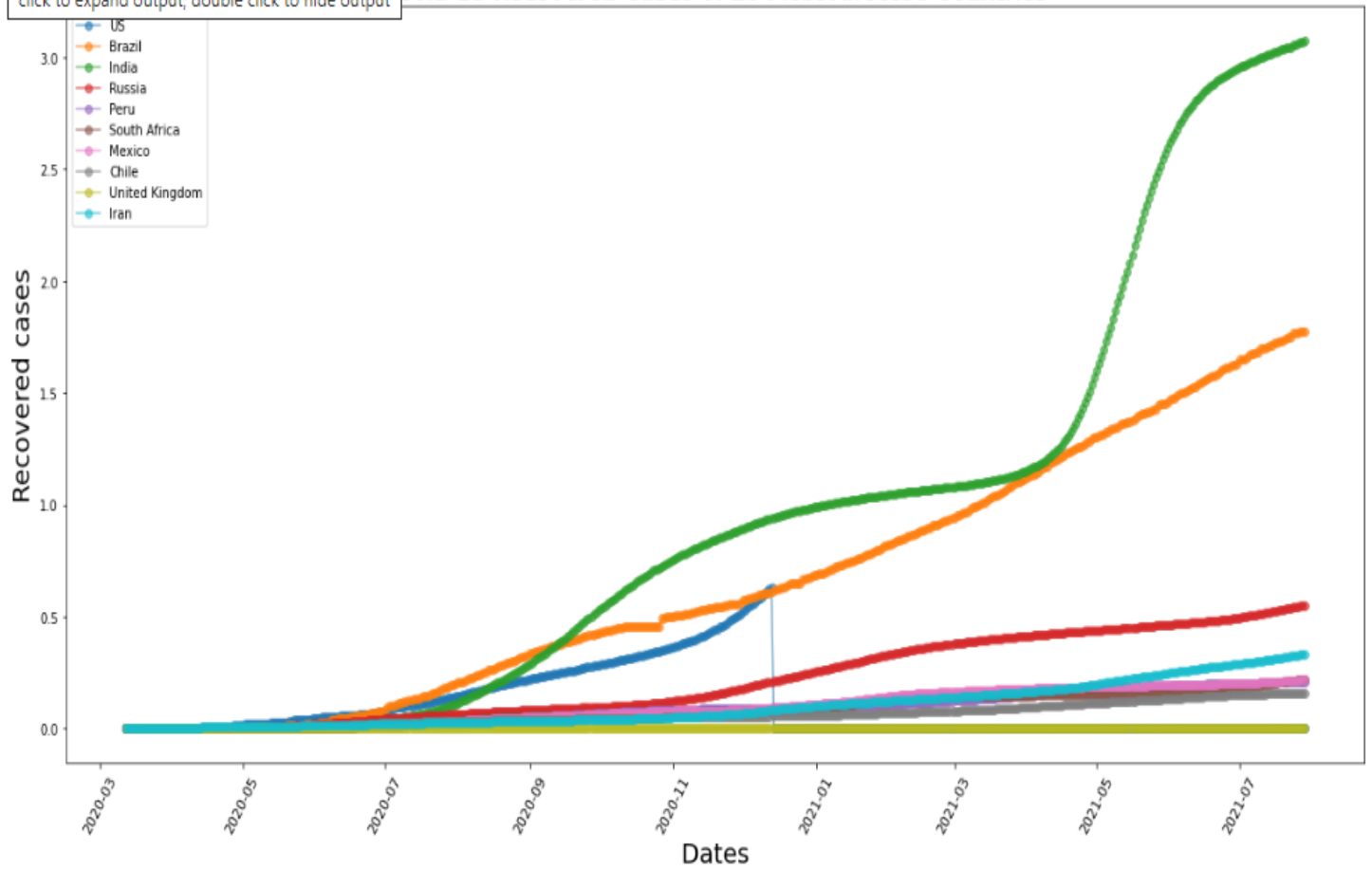
1. Even though the total number of confirmed cases and deaths in the world are monotonically (almost exponentially) increasing, the recovery rate shows some increase whereas the mortality rate shows some decrease.
2. Although US has shown the greatest rise in the number of confirmed cases and deaths.
3. Between 10 most affected countries, India shows the greatest rise in the number of recovered cases, whereas United Kingdom shows very few recoveries.

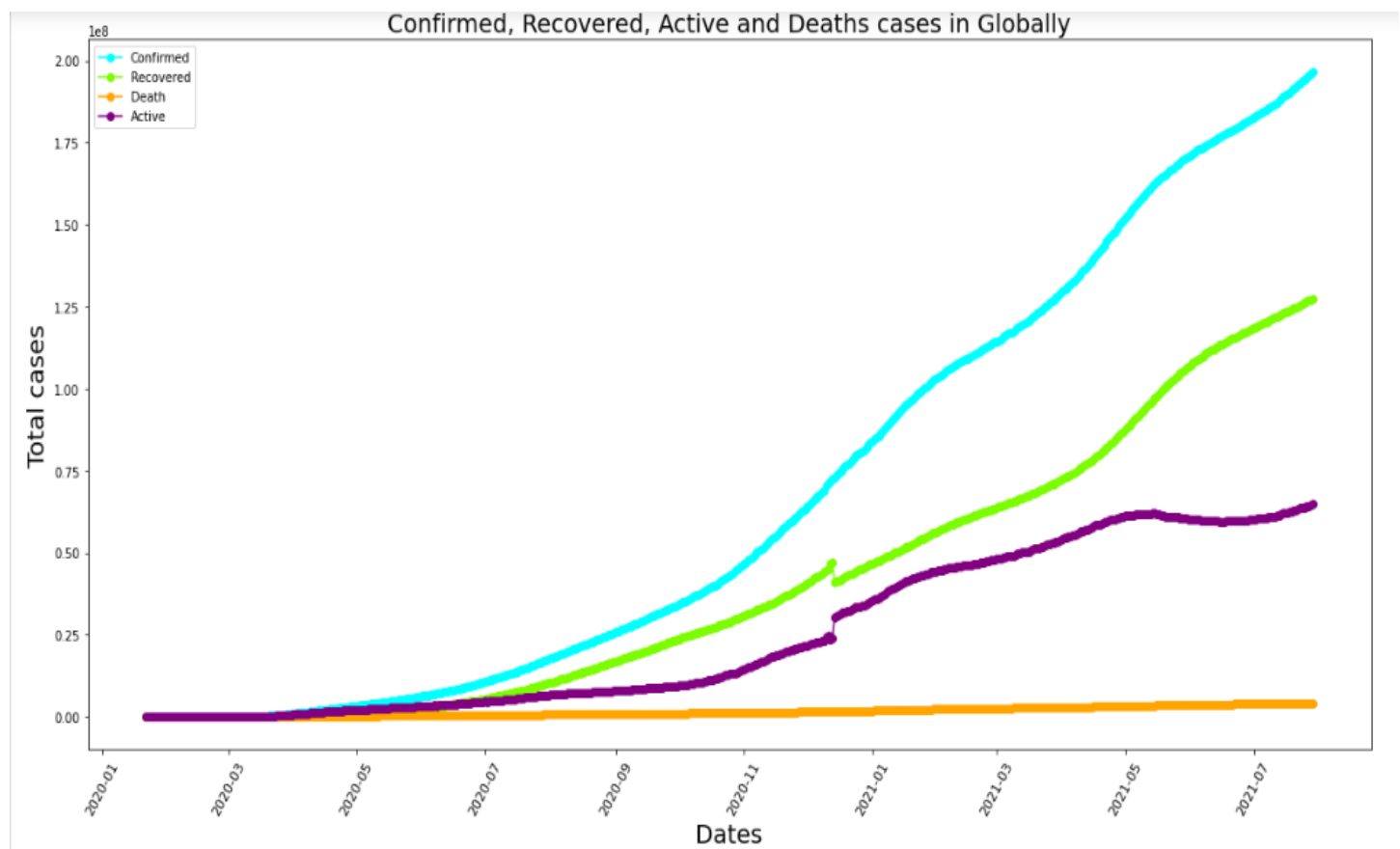
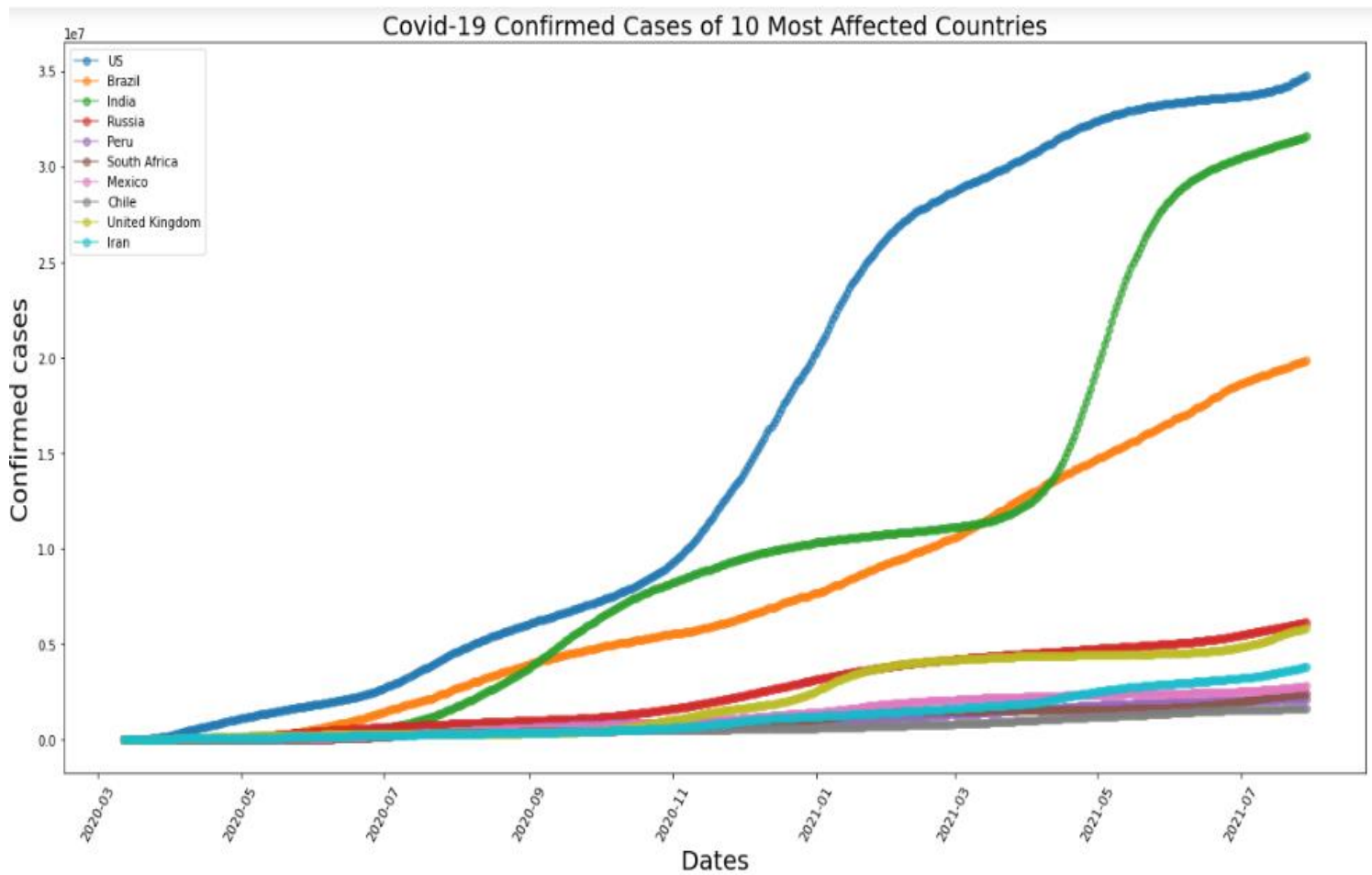
Covid-19 Deaths Cases of 10 Most Affected Countries



click to expand output; double click to hide output

Covid-19 Recovered Cases of 10 Most Affected Countries





click to expand output; double click to hide output

	Country/region	Confirmed	Deaths	Recovered	Active	New confirmed	New deaths	New recovered	Recovery rate(per 100)	Mortality rate(per 100)
0	US	34750860	612122	0	34138738	78170	321	0	0.000000	1.760000
1	India	31572344	423217	30743972	406155	44230	555	42060	-38.660000	1.340000
2	Brazil	19839369	554497	17771228	1513844	42283	1318	28896	89.580000	2.790000
3	France	6142282	111951	412373	5817958	25429	28	576	6.710000	1.820000
4	Russia	6138969	154404	5500409	484156	22720	784	20259	89.600000	2.520000
5	United Kingdom	5828316	129809	20638	5877869	30871	91	914	0.350000	2.230000
6	Turkey	5682630	51184	5443501	187945	22161	60	5463	95.790000	0.900000
7	Argentina	4905925	105113	4542904	257908	14115	291	13645	92.600000	2.140000
8	Colombia	4766829	120126	4557829	88874	9690	325	12483	95.620000	2.520000
9	Spain	4422291	81442	150376	4190473	26689	46	0	3.400000	1.840000
10	Italy	4336906	128029	4130393	78484	6167	19	1825	95.240000	2.950000
11	Iran	3826447	90074	3329065	407308	34433	292	18902	87.000000	2.350000
12	Germany	3772328	91622	3649670	31034	0	-87	-4320	96.750000	2.430000
13	Indonesia	3331206	90552	2686170	554484	43479	1893	45494	80.640000	2.720000
14	Poland	2882630	75257	2653529	153844	165	5	83	92.050000	2.610000
15	Mexico	2810097	239997	2192477	377623	19223	381	11186	78.020000	8.540000
16	South Africa	2422151	71431	2194762	155958	30928	1093	14268	90.610000	2.950000
17	Ukraine	2330440	55489	2253803	21148	1396	40	1085	96.710000	2.380000
18	Peru	2108595	196214	2075361	-162960	722	76	1109	96.420000	9.310000
19	Netherlands	1888741	18102	28149	1842490	4543	4	75	1.490000	0.960000
20	Czechia	1673219	30363	1639849	3007	202	1	151	96.010000	1.810000
21	Chile	1613288	35295	1568294	9699	1371	119	1058	97.210000	2.190000
22	Iraq	1603787	18533	1447160	138094	13259	49	9401	90.230000	1.160000
23	Philippines	1572287	27577	1488437	56273	5620	176	3723	94.670000	1.750000
24	Canada	1437111	26544	0	1410567	912	14	0	0.000000	1.850000
25	Bangladesh	1226253	20255	1050220	155778	15271	239	14336	85.640000	1.650000
26	Belgium	1122951	25235	0	1097716	1862	4	0	0.000000	2.250000
27	Sweden	1099414	14656	0	1084758	619	1	0	0.000000	1.330000
28	Romania	1082880	34275	1047528	1077	170	1	84	96.740000	3.170000
29	Malaysia	1078646	8725	890742	179179	17170	174	12930	82.580000	0.810000

After visualization, we investigated data modeling and prediction based on univariate time series, using Linear regression, Support vector machine, Random forests, XGBoost, Multilayer perceptron (MLP), and a recurrent neural network, Long Short-Term Memory network (LSTM-RNN) to forecast the number of confirmed cases and deaths in the world and some other countries such as India. Some of our results are summarized in the following tables:

Table 1: Prediction errors of total confirmed cases of the world

Regressor	RMSE
Support Vector Machine	68123619.22
Random Forest Regressor	34229279.33
XGBoost	68123619.22

Table 2: Prediction errors of total Deaths of the world

Regressor	RMSE
Support Vector Machine	1289165.65
Linear Regression	1446097.27

Table 3: Accuracy of predicting the total cases of India using MLP and LSTM-RNN

Neural Network	MAPE%	ACCURECY(Percent)
MLP	0.07492	99.99925
LSTM-RNN	0.123195	99.998768

Table 3: Accuracy of predicting the Death cases of India using MLP and LSTM-RNN

Neural Network	MAPE	ACCURECY(Percent)
MLP	0.68879	99.99311
LSTM-RNN	0.780975	99.99219

Table 3: Accuracy of predicting the Confirmed cases of World using MLP and LSTM-RNN

Neural Network	MAPE	ACCURECY(Percent)
MLP	0.29073	99.9970
LSTM-RNN	0.15863	99.9984

Table 3: Accuracy of predicting the Death cases of World using MLP and LSTM-RNN

Neural Network	MAPE	ACCURECY(Percent)
MLP	0.08641	99.99913
LSTM-RNN	0.20931	99.99790

CONFIRMED CASES BY MLP of INDIA

	Confirmed	Confirmed_predicted
2021-07-20	31216337	3.120539e+07
2021-07-21	31257720	3.124217e+07
2021-07-22	31293062	3.127798e+07
2021-07-23	31293062	3.131446e+07
2021-07-24	31371901	3.135093e+07
2021-07-25	31411262	3.138640e+07
2021-07-26	31440951	3.142150e+07
2021-07-27	31484605	3.145737e+07
2021-07-28	31528114	3.149296e+07
2021-07-29	31572344	3.152751e+07
2021-07-30	NaN	3.156227e+07
2021-07-31	NaN	3.159765e+07
2021-08-01	NaN	3.163274e+07
2021-08-02	NaN	3.166792e+07
2021-08-03	NaN	3.170302e+07
2021-08-04	NaN	3.173786e+07
2021-08-05	NaN	3.177266e+07

CONFIRMED CASES BY LSTM of INDIA

	Confirmed	Confirmed_predicted
2021-07-20	31216337	3.123666e+07
2021-07-21	31257720	3.127766e+07
2021-07-22	31293062	3.131917e+07
2021-07-23	31293062	3.136131e+07
2021-07-24	31371901	3.140355e+07
2021-07-25	31411262	3.144584e+07
2021-07-26	31440951	3.148833e+07
2021-07-27	31484605	3.153130e+07
2021-07-28	31528114	3.157446e+07
2021-07-29	31572344	3.161798e+07
2021-07-30	NaN	3.166242e+07
2021-07-31	NaN	3.170573e+07
2021-08-01	NaN	3.174923e+07
2021-08-02	NaN	3.179289e+07
2021-08-03	NaN	3.183672e+07
2021-08-04	NaN	3.188070e+07
2021-08-05	NaN	3.192485e+07

India Death Cases by LSTM

	Deaths	Deaths_predicted
2021-07-20	418480	415310.854507
2021-07-21	418987	415855.996193
2021-07-22	419470	416407.808243
2021-07-23	419470	416968.168241
2021-07-24	420551	417529.714081
2021-07-25	420967	418092.050483
2021-07-26	421382	418656.956211
2021-07-27	422022	419228.285252
2021-07-28	422662	419802.035369
2021-07-29	423217	420380.677125
2021-07-30	NaN	420971.621978
2021-07-31	NaN	421547.447339
2021-08-01	NaN	422125.693795
2021-08-02	NaN	422706.262526
2021-08-03	NaN	423288.955891
2021-08-04	NaN	423873.724461
2021-08-05	NaN	424460.716525
2021-08-06	NaN	423607.749887
2021-08-07	NaN	424024.968812
2021-08-08	NaN	424532.311767
2021-08-09	NaN	425023.744669
2021-08-10	NaN	425479.256425
2021-08-11	NaN	425846.768764
2021-08-12	NaN	426171.887273
2021-08-13	NaN	426534.409212
2021-08-14	NaN	426923.316146
2021-08-15	NaN	427188.500178

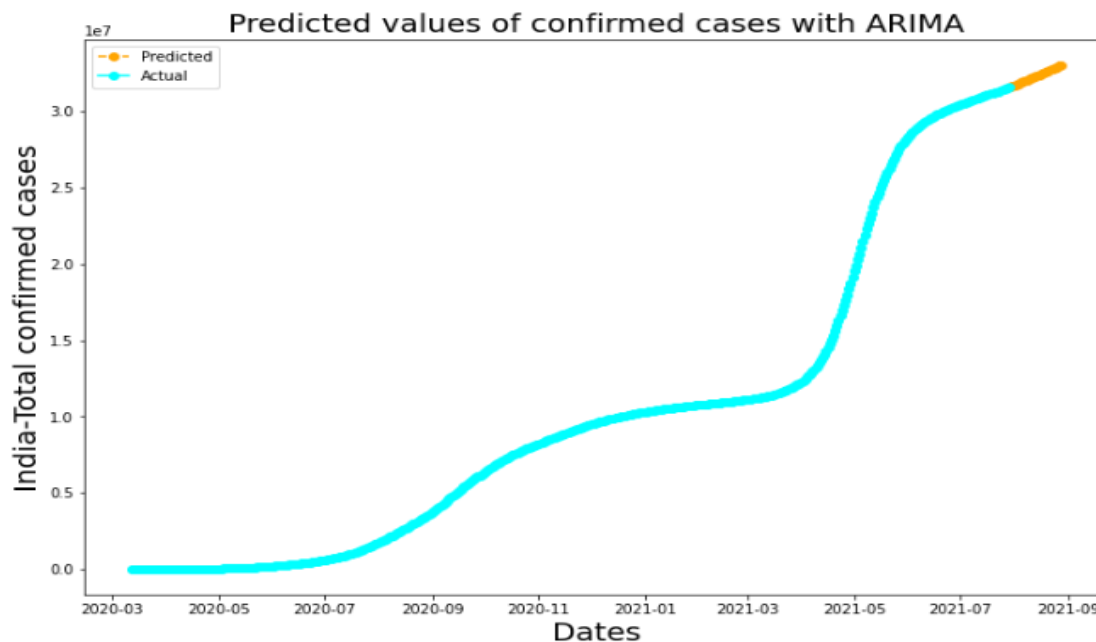
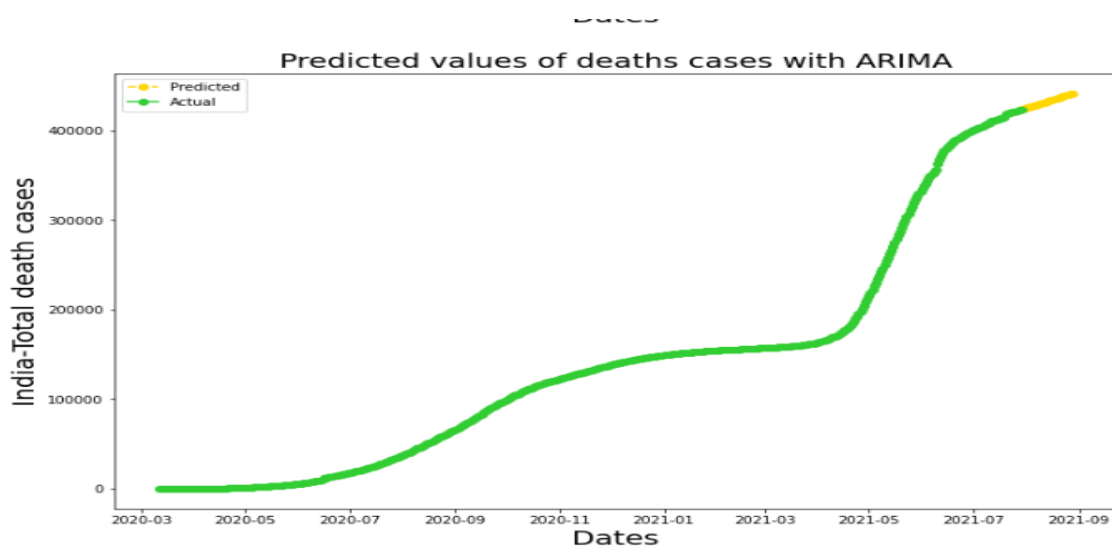
India's Death Cases by MLP

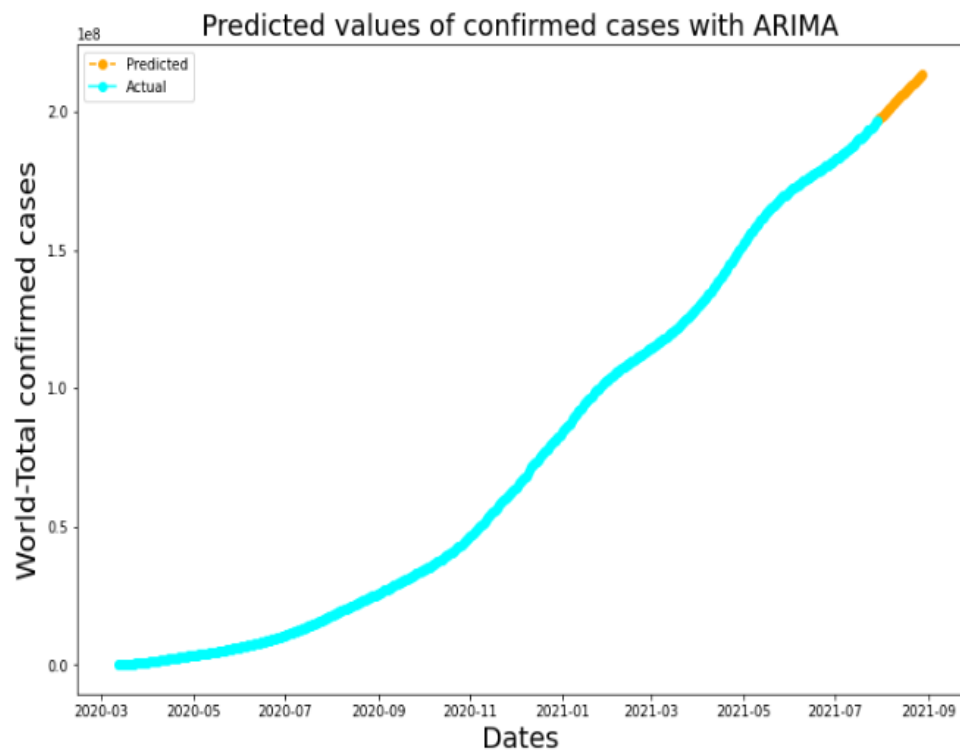
	Deaths	Deaths_predicted
2021-07-20	418480	414877.774939
2021-07-21	418987	415483.789874
2021-07-22	419470	416087.235483
2021-07-23	419470	416630.400765
2021-07-24	420551	417172.923715
2021-07-25	420967	417717.818350
2021-07-26	421382	418266.418743
2021-07-27	422022	418816.847310
2021-07-28	422662	419370.536944
2021-07-29	423217	419929.513459
2021-07-30	NaN	420500.051939
2021-07-31	NaN	421081.411231
2021-08-01	NaN	421664.401057
2021-08-02	NaN	422249.564928
2021-08-03	NaN	422839.867449
2021-08-04	NaN	423435.555672
2021-08-05	NaN	424036.876647

Conclusions and future works:-

As a conclusion based on the analysis of the observations, it seems that even though the total number of confirmed cases and deaths in the world are monotonically (almost exponentially) increasing, the recovery rate shows some increase whereas the mortality rate shows some decrease. On the other hand, by data modeling and prediction based on univariate time series, using Linear regression, Support vector machine, Random forests and XGBoost we concluded that Support vector machine and Random forests performed the best and the worst accuracy, respectively. Moreover, both of Multilayer perceptron and LSTM-RNN performed high accuracy, more than 99.98 in percent.

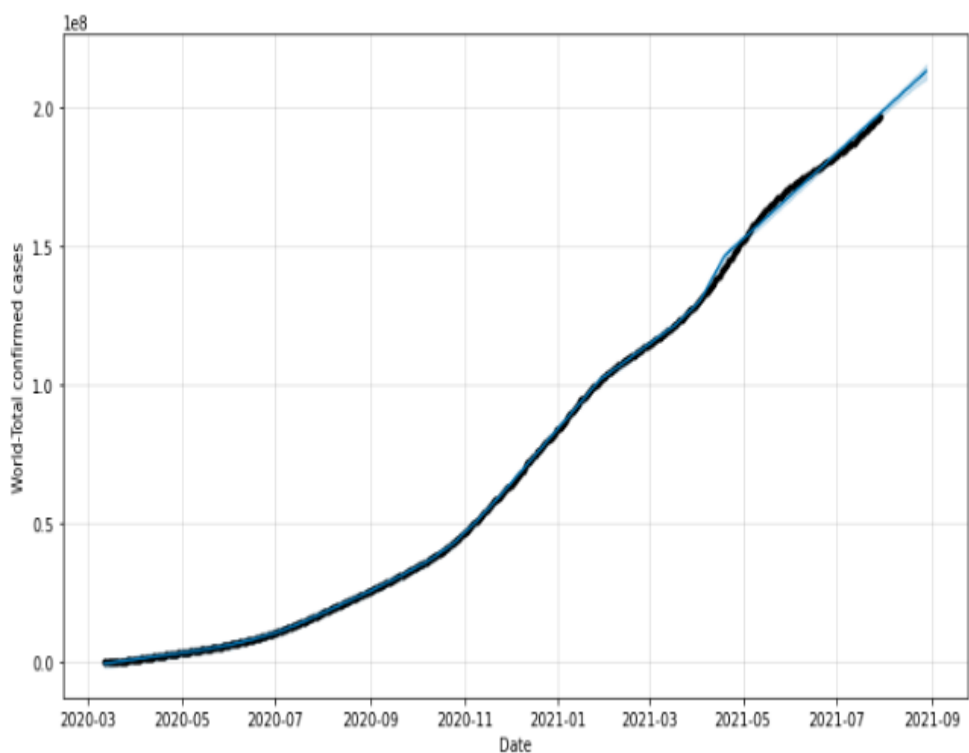
ARIMA & PROPHET FORECASTING:-





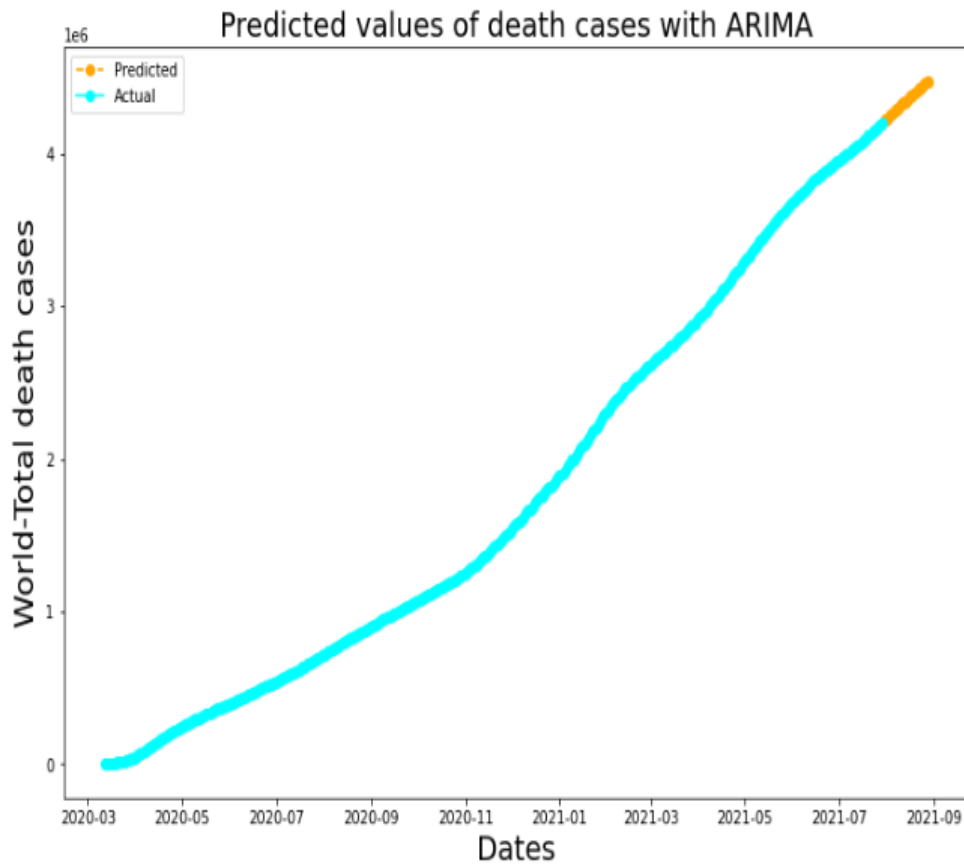
	Date	preidiction
0	2021-07-30	1.972043e+08
1	2021-07-31	1.977626e+08
2	2021-08-01	1.983244e+08
3	2021-08-02	1.989067e+08
4	2021-08-03	1.995150e+08
5	2021-08-04	2.001134e+08
6	2021-08-05	2.006969e+08
7	2021-08-06	2.012665e+08
8	2021-08-07	2.018339e+08
9	2021-08-08	2.024103e+08
10	2021-08-09	2.029914e+08
11	2021-08-10	2.035716e+08
12	2021-08-11	2.041451e+08
13	2021-08-12	2.047123e+08
14	2021-08-13	2.052781e+08
15	2021-08-14	2.058450e+08
16	2021-08-15	2.064132e+08
17	2021-08-16	2.069802e+08
18	2021-08-17	2.075441e+08
19	2021-08-18	2.081051e+08
20	2021-08-19	2.086642e+08

World Confirmed cases forecasted by ARIMA



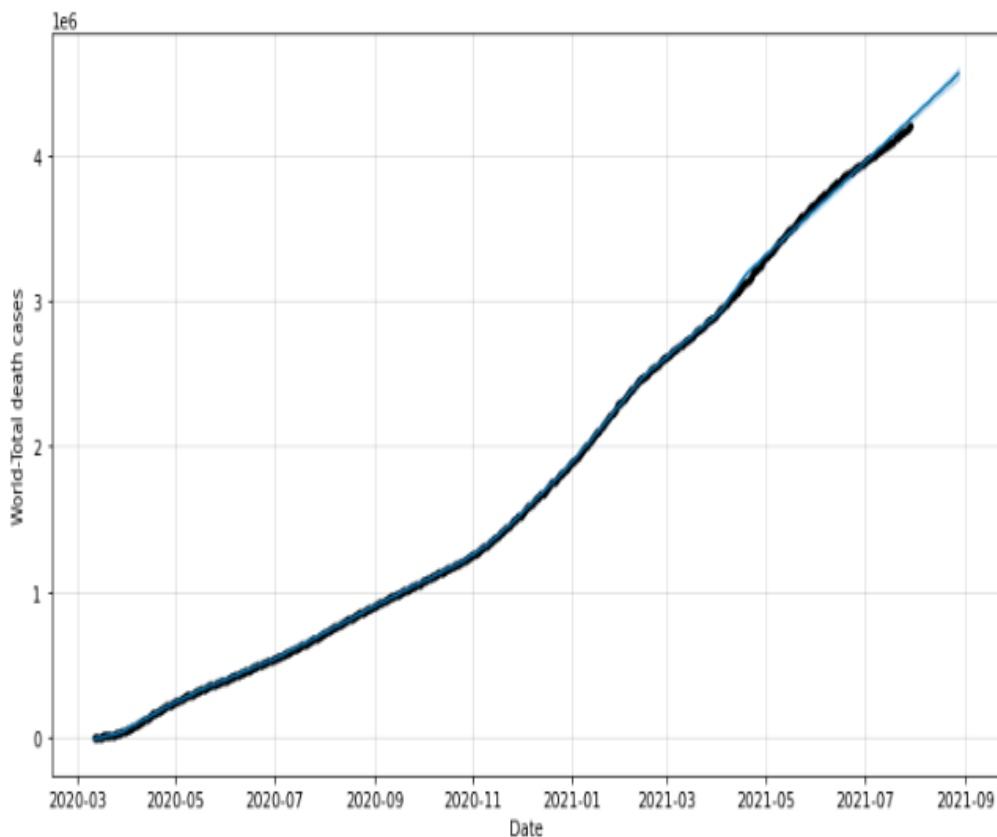
	ds	yhat
505	2021-07-30	1.983193e+08
506	2021-07-31	1.988187e+08
507	2021-08-01	1.992549e+08
508	2021-08-02	1.996984e+08
509	2021-08-03	2.002052e+08
510	2021-08-04	2.007467e+08
511	2021-08-05	2.013106e+08
512	2021-08-06	2.018702e+08
513	2021-08-07	2.023695e+08
514	2021-08-08	2.028058e+08
515	2021-08-09	2.032493e+08
516	2021-08-10	2.037561e+08
517	2021-08-11	2.042976e+08
518	2021-08-12	2.048615e+08
519	2021-08-13	2.054211e+08
520	2021-08-14	2.059204e+08
521	2021-08-15	2.063567e+08
522	2021-08-16	2.068002e+08
523	2021-08-17	2.073070e+08
524	2021-08-18	2.078485e+08
525	2021-08-19	2.084124e+08

World Confirmed Cases Forecasted by Prophet



	Date	prediction
0	2021-07-30	4.209391e+06
1	2021-07-31	4.218767e+06
2	2021-08-01	4.228138e+06
3	2021-08-02	4.237550e+06
4	2021-08-03	4.247410e+06
5	2021-08-04	4.257362e+06
6	2021-08-05	4.267028e+06
7	2021-08-06	4.276465e+06
8	2021-08-07	4.285786e+06
9	2021-08-08	4.295150e+06
10	2021-08-09	4.304585e+06
11	2021-08-10	4.313991e+06
12	2021-08-11	4.323304e+06
13	2021-08-12	4.332522e+06
14	2021-08-13	4.341689e+06
15	2021-08-14	4.350845e+06
16	2021-08-15	4.359989e+06
17	2021-08-16	4.369100e+06
18	2021-08-17	4.378162e+06
19	2021-08-18	4.387176e+06
20	2021-08-19	4.396157e+06

World Deaths Cases by ARIMA

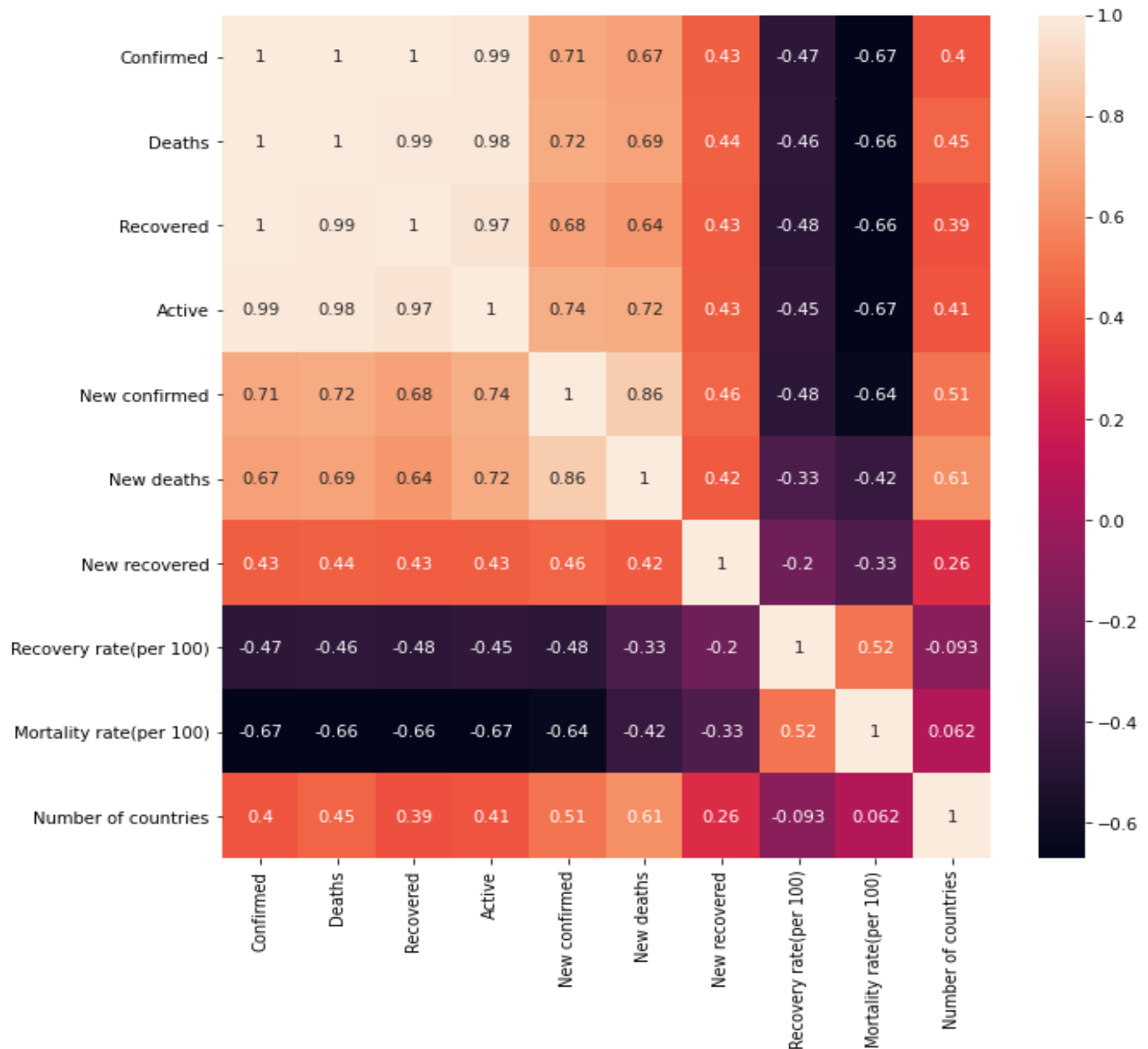


	ds	yhat
505	2021-07-30	4.258170e+06
506	2021-07-31	4.268042e+06
507	2021-08-01	4.275966e+06
508	2021-08-02	4.284816e+06
509	2021-08-03	4.296549e+06
510	2021-08-04	4.308193e+06
511	2021-08-05	4.319730e+06
512	2021-08-06	4.331540e+06
513	2021-08-07	4.341412e+06
514	2021-08-08	4.349336e+06
515	2021-08-09	4.358186e+06
516	2021-08-10	4.369919e+06
517	2021-08-11	4.381563e+06
518	2021-08-12	4.393100e+06
519	2021-08-13	4.404910e+06
520	2021-08-14	4.414782e+06
521	2021-08-15	4.422707e+06
522	2021-08-16	4.431556e+06
523	2021-08-17	4.443289e+06
524	2021-08-18	4.454933e+06
525	2021-08-19	4.466470e+06

World Deaths Cases By Prophet

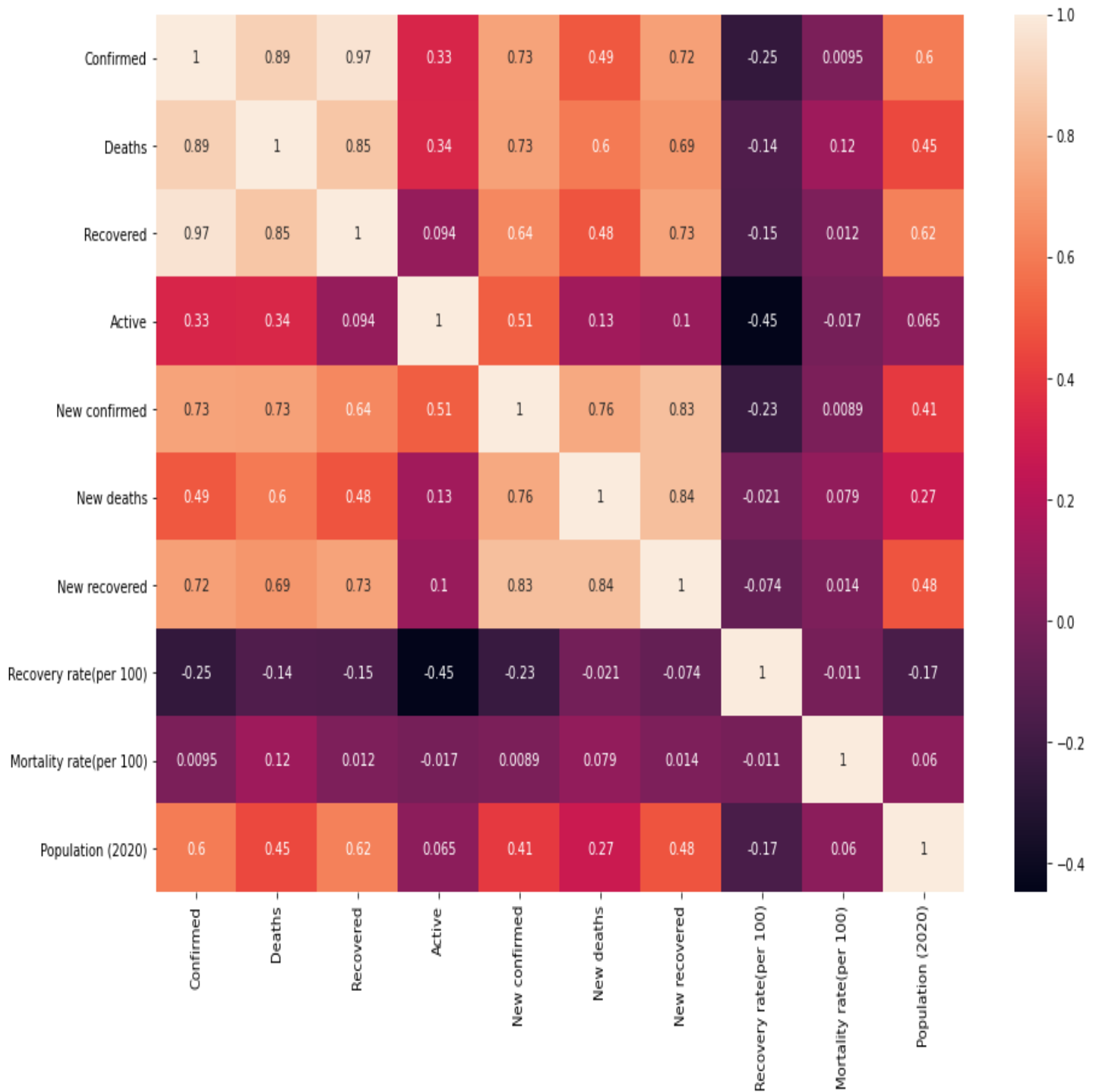
Furthermore, by examining the correlations between the features, it seems that there exist weak correlations between the new parameters, life expectancy, GDP per capita, social support, and freedom to make life choices, generosity, and the primary ones, confirmed, deaths, recovered, active cases, recovery rate and mortality rate. Also, it seems that the correlation between population and confirmed and, population and active cases is moderate (near 0.6). As future works, by considering the population of each country, we may investigate the percentage of total populations that will be affected by COVID-19. Also, the impact of some other parameters in prediction of COVID-19 spread can be considered. Moreover, data modeling and prediction based on multivariate time series using Multilayer perceptron and LSTM-RNN can be considered.

Heat Map for Correlation.

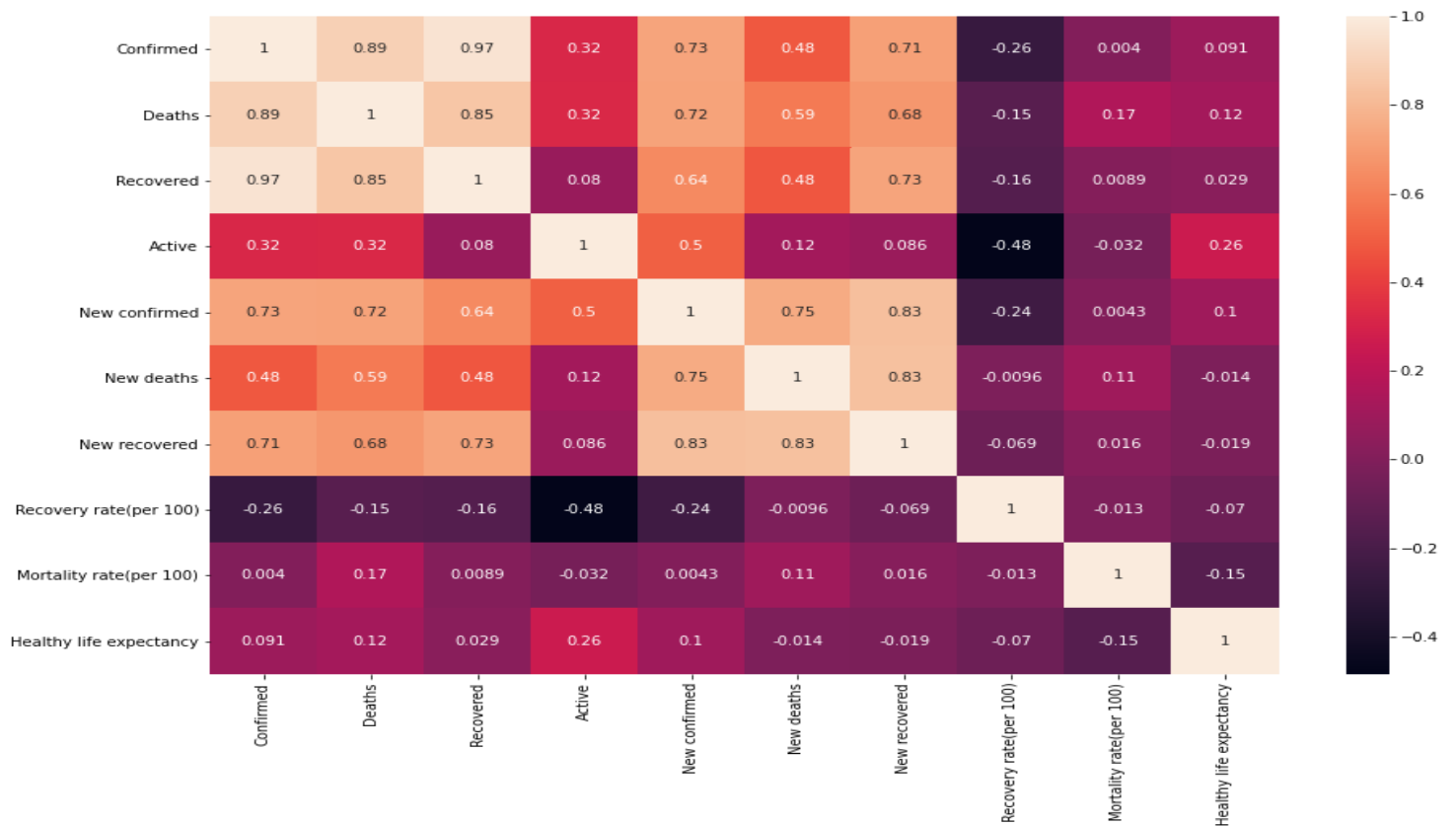
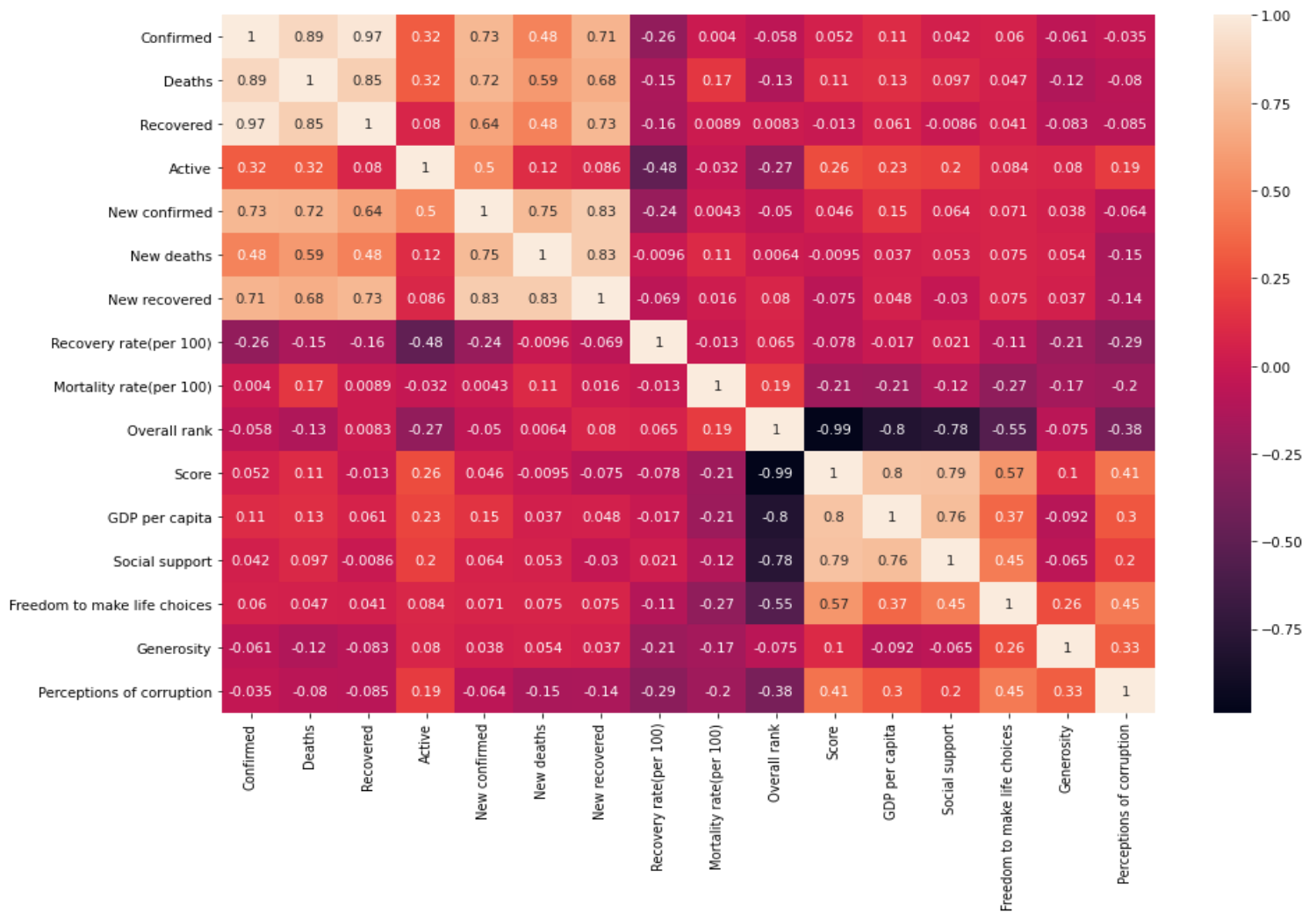


From Here we can understand the relation between different Variables and find out some insight.

Correlaion of Features with Population.



Here we observed that there are moderate correlation between population and confirmed cases as well as deaths cases.



From All heat map we can understand the relation between all variables after that it will help us to take action.

It's a detailed study about my projects. We can conclude that the spread of covid-19 is increasing with respect to time but the good thing is that the mortality rate is decreasing and recovery rate is increasing with respect to time. All this study will help us for taking measures in future. As future works, by considering the population of each country, we may investigate the percentage of total populations that will be affected by COVID-19. Also, the impact of some other parameters in prediction of COVID-19 spread can be considered. Moreover, data modeling and prediction based on multivariate time series using Multilayer perceptron and LSTM-RNN can be considered. And we can also decide how we have to make the strategy for Vaccination drive.

References

- [1] Alon, T. M., et al., The impact of COVID-19 on gender equality, National Bureau of Economic Research, (2020), no. w26947.
- [2] Chen, B., et al., Roles of meteorological conditions in COVID-19 transmission on a worldwide scale, MedRxiv, (2020).
- [3] Fernandes, N., Economic effects of coronavirus outbreak (COVID-19) on the world economy, Available at SSRN 3557504, (2020).
- [4] Fong, S. J., Li, G., and Dey, N., Finding an Accurate Early Forecasting Model from Small Dataset: A Case of 2019-nCoV Novel Coronavirus Outbreak.
- [5] Huang, C., Wang, Y., Li, X., et al., Clinical features of patients infected with 2019 novel coronavirus in Wuhan, China, The Lancet , 395(2020), no. 10223, 497–506.
- [6] Jia, L., Li, K., Jiang, Y., and Guo, X., Prediction and analysis of coronavirus disease 2019, arXiv preprint, (2020).
- [7] Kumar, J., and Hembram, K. P. S. S., Epidemiological study of novel coronavirus (COVID-19), arXiv preprint, (2020).
- [8] World Health Organization (WHO), Naming the coronavirus disease (COVID–19).
- [9] World Health Organization (WHO), Novel Coronavirus–China, Retrieved 9 April 2020.
- [10] Yang, Z., et al., Modified SEIR and AI prediction of the epidemics trend of COVID-19 in China under public health interventions, J Thorac Dis, 12(2020), no. 3, 165.