

COVID-19 Prediction using ARIMA Model

Venkatbharat Polineni

Computer Science and Engineering
CMR University Main Campus
Bangalore, India
bharath.polineni123@gmail.com

Jahnvi K Rao

Information Technology
CMR University Main Campus
Bangalore, India
jahnvi190301@gmail.com

Syed Afshana Hidayathulla

Computer Science and Engineering
CMR University Main Campus
Bangalore, India
syedafshana4812@gmail.com

Abstract—Real-time data has become a dominant aspect for understanding past, present, and future situations. Machine Learning (ML) is one such subject that uses a variety of algorithms to understand the correlation between the given data, visualize the current scenario, and predict the future forecast which is the most crucial part. The entire world is currently undergoing a devastating situation due to the outbreak of novel coronavirus known as COVID-19. The COVID-19 at present has proved that it is a potential threat to human life. To contribute to control the spread and rising number of active cases in India, this study demonstrates the future forecasting of the total number of active cases in India in the upcoming 15 days. The future forecast is predicted using the ARIMA Model (Auto-regressive Integrated Moving Average) with the combination of Facebook Prophet which gives us the highest accuracy. The real-time data collection takes place from various resources after which the data pre-processing and data wrangling takes place. The data set is then split into the training set and testing set. Finally, the model is trained and tested for accuracy. With the completion of testing and training, the model is ready to predict future forecasts. The model also makes note of the predicted and actual values which helps it achieve higher accuracy in the future.

Index Terms—COVID-19, ARIMA Model, FB Prophet, Machine Learning, Time Series Analysis, web scraping, forecasting, R-squared score, Root Mean Squared Error, Mean Squared Error.

I. INTRODUCTION

Machine Learning (ML), considered as one of the most prominent courses in Computer Science in the past few years, has been acclaimed to solve many real-time problems which include image processing, medical diagnosis, financial analysis, etc. Many high profile applications like Autonomous Vehicles (AV), intelligent robots, automatic translations, product recommendations, and climate modelling use ML algorithms that provide them with the highest accuracy. Reinforcement Learning that comes into play while creating ML models, not only avoids the use of traditional step-by-step coding instructions based on logic and if-then rules but also improves their performance over time [1]. Forecasting or predicting the future trend is the best place for ML to showcase its skills [2]. A variety of predictions like weather, the national stock market and many more use ML algorithms to forecast the future so that the necessary action is taken [3]. The speed at which the coronavirus infects humans is rapidly catching up to higher numbers [4]. The need to minimize deaths and stabilizing the country's economy has now become the priority [5]. The virus traverses from a single infected person to a normal human

being through droplets from the mouth or the nose of the infected person, and also when the healthy person comes in contact with the contaminated surface [6]. Safety measures and sanitization must be followed by every individual to overcome this pandemic. The chain must be broken by staying indoors and avoiding visiting places affected by the virus. A frightful situation like this compels for advancements in the field of research and development. Hence, several diligent researchers from all of the science fields have strived to provide any possible solutions [6].

To contribute to the current catastrophe, we have attempted to advance a future forecast for the COVID-19 pandemic. The forecast is regarding the increase in the number of COVID-19 cases for the next 15 days in India. AR (Auto-regression) and MA (Moving Average) models make use of time series analysis. The AR model restores the values from a variable belonging to the early periods as the input for the regression equation which later predicts the output for the upcoming period [7]. The MA model is a time series model that accounts for the possibility of a relationship between a variable and the residuals from the preceding periods. The use of AR and MA models for prediction is not enough due to the lack of accuracy. This gives rise to the ARIMA model which is valid only if the variables are immobile. As of this day, time series performs conversion of the mobile variables into immobile variables using methods such as detrending or differencing for a convincing time series modelling [8]. The conversion now becomes the first initiative to introduce the ARIMA model [9]. ARIMA (p, q, d) denotes the ARMA model with p autoregressive lags, q moving average lags, and the variation in the order of d. The model gets trained with the training set of 85%, the remaining 15% is availed as the testing set, by splitting the data into 85:15 ratio the model gets to train over a substantial amount of data. To find out the performance of the model during the training, we test the model with the unseen test data.

The Facebook Prophet gives the dates and time stamps for the next 15 days. Facebook Prophet (FB Prophet) is an open-source released by Facebook's Core Data Science team which is robust in missing data, outliers, and shifts in the trend. It provides intuitive parameters that are easy to tune, hence we can achieve accurate and fast results [10]. It forecasts time series data based on an additive model where non-linear trends are fit with yearly, weakly, and daily seasonality [11].

The evaluation parameters for the model's performance are R^2 score, MAE (Mean Absolute Error), MSE (Mean Squared Error), and RMSE (Root Mean Squared Error). The useful observations made in this study are listed below:

- The data set is taken from a real-time website which makes the code for this study dynamic.
- Data visualizations are done for a better understanding of the current situation in India.
- ML algorithms need an ample amount of data for better prediction; as the size of the training dataset increases, the model performance increases. Hence for every prediction, the previous day's data is appended to the training data.
- The forecasting based on ML algorithms can be proven very useful in taking protective measures and in guiding the plan of action during pandemics like COVID-19.

II. MATERIALS AND METHODS

A. Dataset

The aim of this research is future forecasting of COVID-19 spread focusing on the total number of active cases in India. The data set used for the study is retrieved from the official government web site [12], which provides the total number of active cases across India for the next 15 days. Further, the data is stored locally in the system with the help of web scraping, an art which extricates a large volume of data from any chosen website, and data can later be stored into the local files in the computer or database. Tables I, II, and III consist of sample datasets.

TABLE I
COVID-19 TIME-SERIES OF INDIA ON DAILY BASIS.

Sl. No.	ds	y
0	2020-01-22	0
1	2020-01-23	0
2	2020-01-24	0
3	2020-01-25	0
.	.	.
149	2020-06-19	395048
150	2020-06-20	410451
151	2020-06-21	425282
152	2020-06-22	440215
153	2020-06-23	456183

TABLE II
COVID-19 DATA ON DAILY BASIS (INDIA).

Date	TC	TD	TR	DC	DD	DR
12-Mar	81	1	4	10	1	0
13-Mar	91	1	10	10	0	6
14-Mar	102	2	10	11	1	0
15-Mar	112	2	13	10	0	3
16-Mar	126	2	14	14	0	1
17-Mar	146	3	15	20	1	1
18-Mar	171	3	15	25	0	0
19-Mar	198	4	20	27	1	5
20-Mar	256	4	23	58	0	3
21-Mar	334	4	23	78	0	0
22-Mar	403	7	23	69	3	0

TABLE III
COVID-19 PATIENT ACTIVE, CURED, AND DEATH ALONG WITH LATITUDE AND LONGITUDE.

State/UT	AC*	C/D/M*	D	TCC*	Date	Lat	Long
AN	11	33	0	44	17-06-2020	11.7401	92.6586
AP	3244	3509	88	6841	17-06-2020	15.9129	79.74
AR	88	7	0	95	17-06-2020	28.218	94.7278
AS	2145	2166	8	4319	17-06-2020	26.2006	92.9376
BR	2093	4644	41	6778	17-06-2020	25.0961	85.3131
CH	50	302	6	358	17-06-2020	30.7333	76.7794
CG	736	1036	9	1781	17-06-2020	21.2787	81.8661
DHND	36	9	0	45	17-06-2020	20.1809	73.0169
DL	26351	16500	1837	44688	17-06-2020	28.7041	77.1025
GA	544	85	0	629	17-06-2020	15.2993	74.124
GJ	5962	17082	1533	24577	17-06-2020	22.2587	71.1924
HR	4406	3748	118	8272	17-06-2020	29.0588	76.0856

B. Time-Series analysis with Auto-regressive Integrated Moving Average (ARIMA)

Auto-regressive Integrated Moving Average is conceivably a category associated with models which unravel a specified census sustained by its previous values, its lags, and also lingered estimate errors. Any 'un-seasonal' statistic which showcases patterns and not belonging to the non-linear noise needs to get shaped based on the present model.

This model is characterized by 3 terms:

p is the order of the AR expression

q is the order of the MA expression

d stating the quantity of differencing needed to form the statistic immobility.

The regression model which makes use of its delay as forecasters are called 'Auto-Regressive'. Regression models recognize which forecasters perform finest when they aren't associated and are independent of each other [13].

Mathematical formula for the AR and MA models:

Auto Regressive (AR only) model is one where Y_t depends only on its lags. That is, Y_t is a function of the 'lags of Y_t '.

$$Y_t = \alpha + \beta_1 Y_{t-1} + \beta_2 Y_{t-2} + \dots + \beta_p Y_{t-p} + \epsilon_1 \quad (1)$$

where, Y_{t-1} is the lag 1 of the series, β_1 is the coefficient of lag 1 that the model estimates, and α is the intercept term, also estimated by the model.

Likewise, a pure Moving Average (MA only) model is one where Y_t depends only on the lagged forecast errors.

$$Y_t = \alpha + \epsilon_t + \phi_1 \epsilon_{t-1} + \phi_2 \epsilon_{t-2} + \dots + \phi_q \epsilon_{t-q} \quad (2)$$

where the error terms are the errors of the auto-regressive models of the respective lags. The errors ϵ_t and ϵ_{t-1} are the errors from the following equations:

$$Y_t = \beta_1 Y_{t-1} + \beta_2 Y_{t-2} + \dots + \beta_0 Y_0 + \epsilon_t \quad (3)$$

$$Y_{t-1} = \beta_1 Y_{t-2} + \beta_2 Y_{t-3} + \dots + \beta_0 Y_0 + \epsilon_{t-1} \quad (4)$$

An ARIMA model is one where the time series was differenced at least once to make it stationary and you combine the AR and the MA terms. So, the equation becomes:

$$Y_t = \alpha + \beta_1 Y_{t-1} + \beta_2 Y_{t-2} + \dots + \beta_p Y_{t-p} \epsilon_t + \phi_1 \epsilon_{t-1} + \phi_2 \epsilon_{t-2} + \dots + \phi_q \epsilon_{t-q} \quad (5)$$

Predicted Y_t = Constant + Linear combination Lags of Y (up to p lags) + Linear Combination of Lagged forecast errors (up to q lags) [13].

To find the order of the ARIMA model using Facebook's Prophet:

The frameworks need to be harmonised concerning the appropriate quandary since the models don't run as intended. An intensive perception of how the fundamental statistic models operate is required to harmonize these systems. The primary input variables to the ARIMA model are the moving average elements, the auto-regressive elements including the utmost degrees of difference. A standard interpreter does not apprehend how to monitor these systems to evade this operation and this sort of system is too laborious to accumulate and estimate [14].

The Prophet Forecasting Model:

A perishable model comprising of three principal elements, namely trend, seasonality, and holidays is practised in the prophet forecasting model [14]. They're consolidated within the subsequent equation:

$$y(t) = g(t) + s(t) + h(t) + \epsilon_t \quad (6)$$

- $g(t)$: section-wise linear or logistic growth curve used to model non-periodic fluctuations in the statistics.
- $s(t)$: periodical variations which can be weekly or yearly seasonality.
- $h(t)$: impacts of holidays provided by the user with variable schedules.
- ϵ_t : error term approximations for any significant changes not implemented by the model [15].

FB prophet attempts to readjust numerous linear and non-linear functions of time by utilizing it as an independent variable. Exponential smoothing, as well as a prophet, practise the same strategy of modelling seasonality as a supplement component. In exponential smoothing the prevailing observations are given more weight in forecasting compared to the earlier observations, since exponentially decreasing weights are ascribed due to the emergence of the observations [16].

C. Evaluation Parameters

For this research, we evaluate the performance of the learning models in terms of R-squared (R^2) score, Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Square error (RMSE).

1) R-Squared Score:

R-squared (R^2) score is a metric of confidence that is easy to compute and intuitive to interpret [17]. It is the degree of how imminent a data point fits the linear regression hence it tells us how good the regression line predicts the real values. It is the coefficient of determination that provides the measure

of variation that has been demonstrated by the self-standing variables in the model. It gives us goodness-of-fit of the model, plus has a score that always prevails between 0 and 1 (0 and 100%). Due to inappropriate fit accompanied by the choice of an erroneous model, the values tend to lie outside the scope of the data. If $R^2 = 0.92$, for instance, a 92% increase in the expense of fuel is due to the increase in the distance traversed. An R^2 score of 1 symbolises that the regression predictions fit the data faultlessly about its mean. A score of 0 indicates that none of the predictions fit the data about its mean. Higher the R^2 score better the model performance. R^2 monotonously improves with the increase in the number of variables but never diminishes. The formula to find the R-squared (R^2) is given by [18]:

$$R^2 = \frac{\text{VarianceExplainedbyModel}}{\text{TotalVariance}} \quad (7)$$

2) Mean Absolute Error (MAE):

MAE could be a quantity accustomed to evaluate how close the projections are to the eventual outcomes, the sum of differences between the model forecasts and the true values [19]. It is an undeviating score hence all the discrete differences are balanced evenly within the mean. The MAE can span from 0 to infinity while it's a negatively-oriented result - the lower the worth, the better the model execution [20].

$$MAE = \frac{1}{n} \sum_{j=1}^n |y_j - \hat{y}_j| \quad (8)$$

3) Mean Squared Error (MSE):

MSE is the mean of the square of the magnitude of the differences between the actual value and the projected value [19]. The mean separation between the precise point and the projection is computed, then squared to urge the error. The squaring is extremely important to get rid of the negative sign, which provides more importance to more substantial variation. This can also be negatively-oriented, that's lesser the value closer to determining the line of the most suitable fit. MSE can be calculated as [21]:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (9)$$

4) Root Mean Square Error (RMSE):

RMSE is the variance or the root of the residuals, where residuals tell how distant the regression curve is from the precise data points [19]. It's the measure of how the residuals disperse around the line of best fit. It will easily be deciphered because its units match the output units. Again, this is usually negatively-oriented and an inferior RMSE value improves the model performance [20].

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (10)$$

III. METHODOLOGY

A study on novel coronavirus formerly known as COVID-19 future forecasting has grabbed special attention from all over the world. To contribute to control the spread and rising number of active cases in India, this study attempts to demonstrate the future forecasting on the total number of active cases in India in the upcoming 15 days. According to the research done for this study, it has come to our notice that the ARIMA model can be the right choice. To start with the data set, they are extracted from a real-time website [12]. The extraction is made possible with the use of web scraping. With the completion of web scraping the data set undergoes data wrangling and data pre-processing after which the data gets stored in the local drive. Now the pre-processed data is visualized for a perfect overview of the data set. The splitting of the data set begins at this moment where it gets partitioned into train data of 85% along with test data of 15%. The 15% test data, taken from the same dataset, is unrevealed to the model during the training period. By hiding a part of the dataset helps in finding out whether the model has over-fit or under-fit, which are few of the biggest complications while training any model. ARIMA Model gets trained by giving in the training data set. After the training, the model is finally ready to undergo the testing phase. Before the model undergoes the testing phase the Facebook Prophet (Facebook prophet is an online open-source for Time-Series Analysis and future forecasting) gives the dates for the next 15 days along with the timestamp to forecast. Finally, the test data set gets appended with the date and time stamp given by Facebook Prophet. As of now, the model is trained on the total number of active cases patterns. ARIMA Model has been evaluated based on important metrics such as R-Squared Score, MAE, MSE, and RMSE and reported in the results. The train data set is uploaded with the daily total number of active cases so that the model trains better every time it predicts the forecast. The sample of the resultant datasets is shown below.

TABLE IV
DATA FOR PREDICTING THE FUTURE.

Sl. No.	ds	yhat	yhat_lower	yhat_upper
155	2020-06-25	396502.678382	389041.104938	407429.501254
156	2020-06-26	404055.577155	393755.847376	415483.032496
157	2020-06-27	412317.438025	398812.575249	42489.019732
158	2020-06-28	420491.433792	406117.647550	432753.741336
159	2020-06-29	427722.229880	413550.636951	442620.591764

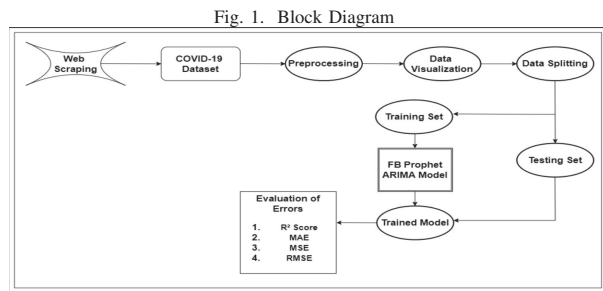
TABLE V
SAMPLE DATA FOR TRAINING THE ARIMA MODEL

SL. No.	Date	Total Confirmed
5	05-Jun	236195
6	06-Jun	246603
7	07-Jun	257485
8	08-Jun	266021
9	09-Jun	276002
10	10-Jun	287158
11	11-Jun	298293
12	12-Jun	309599
13	13-Jun	321638
14	14-Jun	333043
15	15-Jun	343075

TABLE VI
DATES FOR FUTURE PREDICTION.

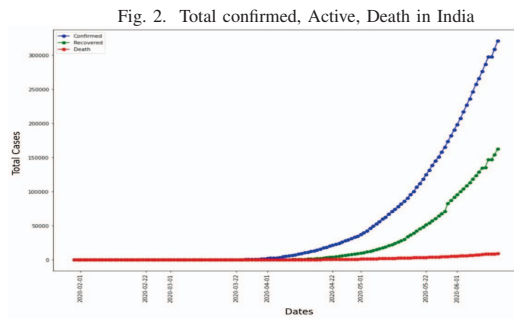
Sl. No.	ds
145	2020-06-15
146	2020-06-16
147	2020-06-17
148	2020-06-18
149	2020-06-19
150	2020-06-20
151	2020-06-21
152	2020-06-22
153	2020-06-23
154	2020-06-24
155	2020-06-25
156	2020-06-26
157	2020-06-27
158	2020-06-28
159	2020-06-29

The block diagram below depicts the process of future forecasting in detail (Figure 1).

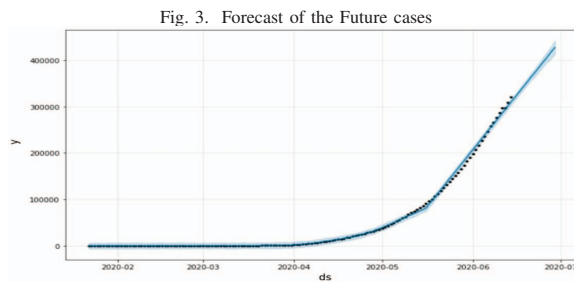


IV. RESULTS AND DISCUSSION

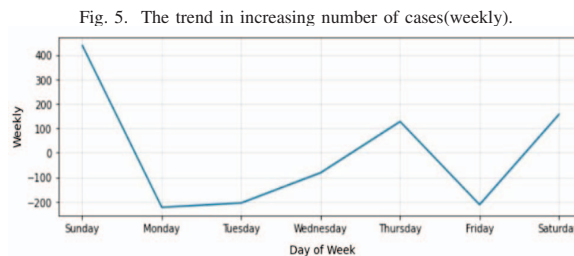
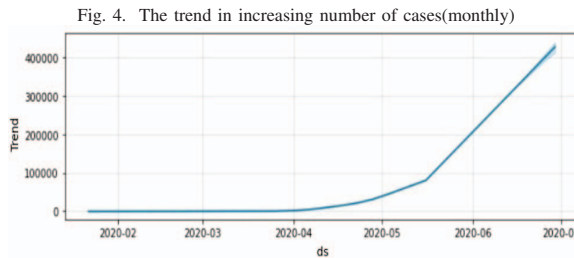
Prognosticating the total number of confirmed cases concerning COVID-19 in the upcoming 15 days is the main objective of this study. The data set consists of the total number of confirmed, active and death cases as shown in Figure-2.



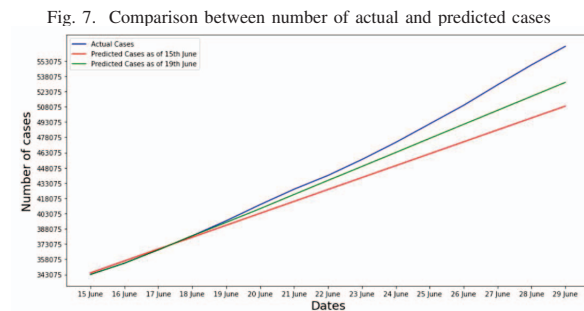
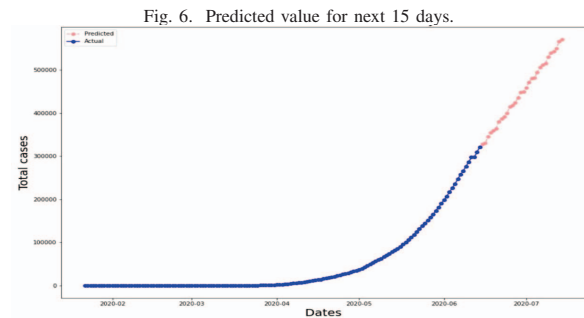
As per the objective of this study, only the required information gets retained in the data set. Collecting the required amount of information for almost six months now, the trend can be observed with the assistance of time series analysis as shown below in Figure-3.



The trend in the rising number of active cases noticed monthly and weekly make the ARIMA model predict the future forecast.



A. Total number of confirmed cases for future forecasting
The study projects the predictions on the number of confirmed cases. According to the R^2 score, the ARIMA model performed best for this scenario which is shown below in Figure-6.



B. Effectively evaluating the model:
The prediction made by the ARIMA model for the upcoming 15 days gets noted down in a data set which further compares itself with the actual number of active cases in India, the predicted cases and the confirmed cases undergo cross-checking after 15 days from the day of prediction. The accuracy results as of 15th June and as of 19th June are exposed below. The results exhibited below clearly depicts that the predictions made as of 19th June have higher accuracy because the model has rectified the errors from the previous forecast. After learning the blunders, the model trains on the mistakes to improve its performance by tuning the trends.

1) Predictions done on 15th June:

TABLE VII
ACCURACY AS OF 29TH JUNE.

Model	R^2	MAE	MSE	RMSE
ARIMA	0.84	18054.95	619782209.7	24895.42

TABLE VIII
ACCURACY AS OF 30TH JUNE.

Model	R^2	MAE	MSE	RMSE
ARIMA	0.83	20912.2	825828637.3	28737.23

2) Predictions done on 19th June:

TABLE IX
ACCURACY AS OF 29TH JUNE.

Model	R^2	MAE	MSE	RMSE
ARIMA	0.96	5588.87	66378229.3	8147.28

TABLE X
ACCURACY AS OF 30TH JUNE.

Model	R^2	MAE	MSE	RMSE
ARIMA	0.96	7084.05	101932387	10096.15

V. CONCLUSION

The COVID-19 continues the potential threat across the planet because the number of victims and deaths rapidly rises. The COVID-19 has clearly shown its role within the recent decline of the country's economy on an enormous scale. It has the potential to infect anybody. There are grave concerns that the economic fallout from COVID-19 may be comparable to that of the great depression [22]. This study predicts the overall number of active cases for the next 15 days across India using the ML approach. The result obtained from this study depicts that the ARIMA model is best fitted to this scenario. The prediction of the model to the present situation will be helpful to know the upcoming situation. Overall, this study can help the authorities to acquire caution which might help arrange to contain the COVID-19 crisis. This study is enhanced continuously within the future course. The model could be further incorporated to predict the economy of the country during a pandemic.

REFERENCES

- [1] S. Makridakis, E. Spiliotis, and V. Assimakopoulos, "Statistical and machine learning forecasting methods: Concerns and ways forward," *PLoS one*, vol. 13, no. 3, p. e0194889, 2018.
- [2] J.-H. Han and S.-Y. Chi, "Consideration of manufacturing data to apply machine learning methods for predictive manufacturing," in *2016 Eighth International Conference on Ubiquitous and Future Networks (ICUFN)*. IEEE, 2016, pp. 109–113.
- [3] F. E. Harrell Jr, K. L. Lee, D. B. Matchar, and T. A. Reichert, "Regression models for prognostic prediction: advantages, problems, and suggested solutions," *Cancer treatment reports*, vol. 69, no. 10, pp. 1071–1077, 1985.
- [4] A. Remuzzi and G. Remuzzi, "Covid-19 and italy: what next?" *The Lancet*, 2020.
- [5] R. M. Anderson, H. Heesterbeek, D. Klinkenberg, and T. D. Hollingsworth, "How will country-based mitigation measures influence the course of the covid-19 epidemic?" *The Lancet*, vol. 395, no. 10228, pp. 931–934, 2020.
- [6] "Predicting the time period of extension of lockdown due to increase in rate of covid-19 cases in india using machine learning - sciencedirect," <https://www.sciencedirect.com/science/article/pii/S2214785320363914>, (Accessed on 12/15/2020).
- [7] A. Meyler, G. Kenny, and T. Quinn, "Forecasting irish inflation using arima models," 1998.
- [8] G. P. Zhang, "Time series forecasting using a hybrid arima and neural network model," *Neurocomputing*, vol. 50, pp. 159–175, 2003.
- [9] C. Javier, E. Rosario, J. Francisco, and J. C. Antonio, "Arima models to predict next electricity price," *IEEE Transactions on Power Systems*, vol. 18, no. 3, pp. 1014–1020, 2003.
- [10] S. J. Taylor and B. Letham, "Forecasting at scale," *The American Statistician*, vol. 72, no. 1, pp. 37–45, 2018.
- [11] "Prophet — prophet is a forecasting procedure implemented in r and python. it is fast and provides completely automated forecasts that can be tuned by hand by data scientists and analysts." <https://facebook.github.io/prophet/>.
- [12] "Mohfw — home," <https://www.mohfw.gov.in/>, (Accessed on 12/15/2020).
- [13] "Arima model - complete guide to time series forecasting in python — ml+," <https://www.machinelearningplus.com/time-series/arima-model-time-series-forecasting-python/>.
- [14] I. Yenidogan, A. Cayir, O. Kozan, T. Dag, and C. Arslan, "Bitcoin forecasting using arima and prophet," in *2018 3rd International Conference on Computer Science and Engineering (UBMK)*. IEEE, 2018, pp. 621–624.
- [15] S. Taylor and B. Letham, "Forecasting at scale. peerj preprints 5: e3190v2 (2017)."
- [16] "Time series forecasts using facebook's prophet," <https://www.analyticsvidhya.com/blog/2018/05/generate-accurate-forecasts-facebook-prophet-python-r/#~:text=Prophet%20is%20an%20open%20source,of%20custom%20seasonality%20and%20holidays!>
- [17] J. Lupón, H. K. Gaggin, M. De Antonio, M. Domingo, A. Galán, E. Zamora, J. Vila, J. Peñafiel, A. Urrutia, E. Ferrer *et al.*, "Biomarker-assist score for reverse remodeling prediction in heart failure: the st2-r2 score," *International journal of cardiology*, vol. 184, pp. 337–343, 2015.
- [18] O. Renaud and M.-P. Victoria-Feser, "A robust coefficient of determination for regression," *Journal of Statistical Planning and Inference*, vol. 140, no. 7, pp. 1852–1862, 2010.
- [19] C. J. Willmott and K. Matsuura, "Advantages of the mean absolute error (mae) over the root mean square error (rmse) in assessing average model performance," *Climate research*, vol. 30, no. 1, pp. 79–82, 2005.
- [20] W. Wang and Y. Lu, "Analysis of the mean absolute error (mae) and the root mean square error (rmse) in assessing rounding model," in *IOP Conference Series: Materials Science and Engineering*, vol. 324, no. 1, 2018, p. 012049.
- [21] H. Witzgall and J. Goldstein, "A realizable mean square error estimator applied to rank selection," in *Conference Record of the Thirty-Sixth Asilomar Conference on Signals, Systems and Computers, 2002.*, vol. 1. IEEE, 2002, pp. 881–884.
- [22] "Global impact of new corona virus and population issues — inter press service," <http://www.ipsnews.net/2020/05/global-impact-new-corona-virus-population-issues/>.