



Modeling and Forecasting of COVID-19 Growth Curve in India

Vikas Kumar Sharma¹ · Unnati Nigam¹

Received: 11 June 2020 / Accepted: 24 August 2020 / Published online: 5 September 2020
© Indian National Academy of Engineering 2020

Abstract

In this article, we analyze the growth pattern of COVID-19 pandemic in India from March 4 to July 11 using regression analysis (exponential and polynomial), auto-regressive integrated moving averages (ARIMA) model as well as exponential smoothing and Holt–Winters models. We found that the growth of COVID-19 cases follows a power regime of (t^2, t, \dots) after the exponential growth. We found the optimal change points from where the COVID-19 cases shifted their course of growth from exponential to quadratic and then from quadratic to linear. After that, we saw a sudden spike in the course of the spread of COVID-19 and the growth moved from linear to quadratic and then to quartic, which is alarming. We have also found the best fitted regression models using the various criteria, such as significant p-values, coefficients of determination and ANOVA, etc. Further, we search the best-fitting ARIMA model for the data using the AIC (Akaike Information Criterion) and provide the forecast of COVID-19 cases for future days. We also use usual exponential smoothing and Holt–Winters models for forecasting purpose. We further found that the ARIMA (5, 2, 5) model is the best-fitting model for COVID-19 cases in India.

Keywords COVID-19 · Regression analysis · Exponential growth · Polynomial growth · ANOVA · ARIMA · Exponential smoothing and holt–winters models · Prediction · Forecast

Introduction

The COVID-19 pandemic has created a lot of havoc in the world. It is caused by a virus called SARS-CoV-2, which comes from the family of coronaviruses and is believed to be originated from the unhygienic wet seafood market in Wuhan, China but it has now infected around 215 countries of the world. With more than 13.2 million people affected around the world and more than 575,000 deaths (As of July 14, 2020), it has forced people to stay in their homes and has caused huge devastation in the world economy (Singh and Singh 2020; Ministry of Health and Family Welfare 2020; Gupta et al. 2019).

In India, the first case of COVID-19 was reported on 30th January, which was linked to the Wuhan city of China (as the patient has travel history to the city). On 4th March, India saw a sudden hike in the number of cases and since then, the numbers are increasing day by day. As of 14th July, India

has more than 908,000 cases with more than 23,000 deaths and is world's 3rd most infected country (<https://www.worldometers.info/coronavirus/>).

Since the outbreak of the pandemic, scientists across the world have been indulged in the studies regarding the spread of the virus. Lin et al. (2020) suggested the use of the SEIR (Susceptible–Exposed–Infectious–Removed) model for the spread in China and studied the importance of government-implemented restrictions on containing the infection. As the disease grew further, Ivorra et al. (2019) suggested a θ -SEIHRD model that took into account various special features of the disease. It also included asymptomatic cases into account (around 51%) to forecast the total cases in China (around 168,500). Giordano et al. (2003) also suggested an extended SIR model called SIDHARTHE model for cases in Italy which was more customized for COVID-19 to effectively model the course of the pandemic to help plan a better control strategy.

Petropoulos and Makridakis (2020) suggested the use of exponential smoothing method to model the trend of the virus, globally. Kumar et al. (2020) gave a review on the various aspects of modern technology used to fight against COVID-19 crisis.

✉ Vikas Kumar Sharma
vikasstats@rediffmail.com

¹ Department of Statistics, Institute of Science, Banaras Hindu University, Varanasi, India

Apart from the epidemiological models, various data-oriented models were also suggested to model the cases and predict future cases for various disease outbreaks from time to time. Various time-series models were also suggested to model the cases and predict future cases. ARIMA and Seasonal ARIMA models are widely used by researchers to model and predict the cases of various outbreaks. In 2005, Earnest et al. (2005) conducted a research to model and predict the cases of SARS in Singapore and predict the hospital supplies needed using this model. Gaudart et al. (2009) modelled malaria incidence in the Savannah area of Mali using ARIMA. Zhang et al. (2013) compared Seasonal ARIMA model with three other time-series models to compare Typhoid fever incidence in China. Polwiang (2020) also used this model to determine the time-series pattern of Dengue fever in Bangkok.

For COVID-19 as well, various researchers tried to model the cases through ARIMA. Ceylan (2020) suggested the use of Auto-Regressive Integrated Moving Average (ARIMA) model to develop and predict the epidemiological trend of COVID-19 for better allocation of resources and proper containment of the virus in Italy, Spain and France. Chintalapudi (2020) suggested its use for predicting the number of cases and deaths post 60-days lockdown in Italy. Fanelli and Francesco (2020) analyzed the dynamics of COVID-19 in China, Italy and France using iterative time-lag maps. It further used SIRD model to model and predict the cases and deaths in these countries. Zhang et al. (2020) developed a segmented Poisson model to analyze the daily new cases of six countries to find a peak point in the cases.

Since the spread of the virus started to grow in India, various measures were taken by the Indian Government to contain it. A nationwide lockdown was announced on March 25 to April 14, which was later extended to May 3. The whole country was divided into containment zones (where large number of cases were observed from a relatively smaller region), red zones (districts where risk of transmission was high and had higher doubling rates), green zones (districts with no confirmed case from last 21 days) and orange zones (which did not fall into the above three zones). After the further extension of the lockdown till May 17, various economic activities were allowed to start (with high surveillance) in areas of less transmission. Further, the lockdown was extended to May 31 and some more economic activities have been allowed as per the transmission rates, which are the rates at which infectious cases cause new cases in the population, i.e. the rate of spread of the disease. This was further extended to June 8, with very less rules and especially the states were given the responsibility of setting the lockdown rules. The air and rail transport became open for general public. Post June 8, we see that the restrictions are nominal with even shopping malls and religious places open for general public. Now, the responsibility of imposing restrictions lies with the respective State Governments.

On the other hand, Indian scientists and researchers are also working on addressing the issues arising from the pandemic, including production of PPE kits and test kits as well as studying the behaviour of spread of the disease and other aspects of management. Various mathematical and statistical methods have been used for predicting the possible spread of COVID-19. The classical epidemiological models (SIR, SEIR, SIQR etc.) suggested the increasing trend of the virus and predicted the peaks of the pandemic. Early researches showed the pandemic to reach its peak by mid-May. They also showed that the basic reproduction number (R_0) and the doubling rates are lower in India, with comparison to European nations and the USA. A tree-based model was proposed by Arti and Bhatnagar (2020) and Bhatnagar (2020) to study and predict the trends. They suggest that lockdown and social distancing in India have played a significant role to control the infection rates. But now, as the lockdown restrictions are minimal, the cases in India are growing at an alarming high rate. Chatterjee et al. (2020) suggest growth of the pandemic through power law and its saturation at the later stages. Due to the complexities in the epidemic models of COVID-19, various researchers have been focusing on the data to forecast the future cases. Chatterjee et al. (2020), Verma et al. (2020) and Ziff and Ziff (2020) suggest that after exponential growth, the total count follows a power regime of t^3 , t^2 , t and \sqrt{t} before flattening out, where ' t ' refers to time. It can, therefore, be realized that there is an urgent need to model and forecast the growth of COVID-19 in India as the virus is in the growing stage here.

In India, the most affected states are Maharashtra with over 260,000 cases (as of 14 July 2020), Tamil Nadu (around 142,000 cases), Delhi (around 113,000 cases) and Gujarat (around 42,000 cases). The greatest number of cases per million has been seen in the national capital of Delhi (5740 cases per million) (Refer https://nhm.gov.in/New_Updates_2018/Report_Population_Projection_2019.pdf for population estimates). Many states and union territories like, Kerala, Karnataka, Andaman and Nicobar Islands, Daman and Diu, etc. which had recovered from majority of the cases have experienced a second wave of infections. This might be attributed to decreased travel restrictions and minimal lockdown measures. In their research, Singh and Jadaun (2020) studied the significance of lockdown in India and suggested that the new COVID-19 cases would stop by the end of August in India with around 350,000 total cases. While some states may see an early stopping of new cases, such as Telangana (mid-June), Uttar Pradesh and West Bengal (July end) etc., the badly affected states of Maharashtra, Tamil Nadu and Gujarat will achieve this by August end.

Since a proven vaccine and medication is yet to be developed by the researchers then in such a scenario, modelling the present situation and forecasting the future outcome becomes crucially important to utilize our resources in the

most optimal way. Therefore, the article aims to study the growth curve of COVID-19 cases in India and forecast its future course. Since the disease is still in its growing age and very dynamic in nature, no model can guarantee for perfect validity for future. We, therefore, need to develop the understanding of the present situation of the pandemic.

In this article, we first study the growth curve using regression methods (exponential, linear and polynomial etc.) and propose an optimal model for fitting the cases till July 10. Further, we propose the use of time-series models for forecasting the future observations on COVID-19 cases. Here, we reach the best-fitted ARIMA model for forecasting the COVID-19 cases. We also compare these results with Exponential Smoothing (Holt–Winters) model. This study will help us to understand the course of spread of SARS-CoV-2 in India better and help the government and the people to optimally use the resources available to them.

Statistical Methodologies

In this section, we briefly present the statistical techniques used for analyzing the COVID-19 cases in India. Here, we used usual regression (exponential, polynomial), times-series (ARIMA) and exponential smoothing models.

Exponential–Polynomial Regression

Regression is a statistical technique that attempts to estimate the strength and nature of relationship between a dependent variable and a series of independent variables. Regression analyses may be linear and non-linear. A regression is called linear when it is linear in parameters, e.g. $y = \beta_0 + \beta_1 t + \epsilon$ and $y = \beta_0 + \beta_1 t + \beta_2 t^2 + \beta_3 t^3 + \epsilon$, $\epsilon \sim N(0, \sigma^2)$, where y is response variable, t denotes the independent variable, β_0 is the intercept and other β s are known as slopes.

A non-linear regression is a regression when it is non-linear in its parameters, e.g. $y = \theta_1 e^{\theta_2 x} + \epsilon$. In the beginning of the spread of a disease, we see that the new cases are directly proportional to the existing infected cases and may be represented by $\frac{dy(t)}{dt} = ky(t)$, where k is the proportionality constant. Solving this differential equation, we get that, at the beginning of a pandemic,

$$y(t) = Ae^{kt}.$$

Thus, at the beginning of a disease, the growth curve of the cases grows exponentially.

As the disease spreads in a region, governments start to take action and people start becoming conscious about the disease. Thus, after some time, the disease starts to follow a polynomial growth rather than continuing to grow exponentially.

In order to fit an exponential regression to our data, we linearize the equation by taking the natural logarithm of the equation and convert it to a linear regression in first order.

We estimate the parameters of a linear regression of order p as follows:

Let the model of linear regression of order p be: $y_i = \beta_0 + \sum_{j=1}^p \beta_j x_i^j + \epsilon$ with $\epsilon \sim N(0, \sigma^2)$ and $i = 1, 2, \dots, N$. Let $E = \sum_{i=1}^N \left\{ y_i - \beta_0 - \sum_{j=1}^p \beta_j x_i^j \right\}^2$ represent the residual sum of square (RSS).

We get the best estimates of these coefficients by solving the following normal equations: $\frac{\partial E}{\partial \beta_0} = 0, \frac{\partial E}{\partial \beta_1} = 0, \dots, \frac{\partial E}{\partial \beta_p} = 0$, which minimizes RSS. This technique is referred to as the ordinary least squares (OLS). We will use this technique of the OLS to estimate the coefficients of our proposed model. (Refer Montgomery et al. (2012)).

Since we know that the growth curve of the disease changes after some time point, exponential to polynomial, we propose to use the following joint regression model with change point μ ,

$$y = \begin{cases} f_1(t); & t \leq \mu, \\ f_2(t); & t > \mu, \end{cases} \tag{1}$$

where we take $f_1(t) = \theta_1 e^{\theta_2 t}$, $f_2(t) = \beta_0 + \beta_1 t + \beta_2 t^2 + \dots + \beta_p t^p + \epsilon$, $\epsilon \sim N(0, \sigma^2)$ and p is the order of the polynomial regression model and t stands for the time (an independent variable).

During the analysis, we found that a suitable choice of $f_2(t)$ is a quadratic or a cubic model. Once the order of the polynomial is kept fixed, an optimum value of the change point can be obtained by minimizing the residuals/errors. We can obtain the OLS estimates of the parameters of the model (1) as given below:

The least square estimates (LSEs) of the parameters, $\Theta = \{\theta_1, \theta_2, \mu, \beta_0, \beta_1, \beta_2, \beta_3, \dots, \beta_p\}$ can be obtained by minimizing the residual sum of squares (RSS) as given by:

$$RSS(\Theta) = \sum_{i=1}^{\mu} (y_i - \hat{y}_i^{exp})^2 + \sum_{i=\mu+1}^N (y_i - \hat{y}_i^{poly})^2, \tag{2}$$

where \hat{y}_i^{exp} and \hat{y}_i^{poly} are the estimates value of y_i from the exponential and polynomial regression models, respectively, and N is the size of the dataset.

The LSEs of $\Theta = \{\theta_1, \theta_2, \mu, \beta_0, \beta_1, \beta_2, \beta_3, \dots, \beta_p\}$ can be obtained as the simultaneous solution of the following normal equations, $\frac{\partial RSS(\Theta)}{\partial \theta_1} = 0, \frac{\partial RSS(\Theta)}{\partial \theta_2} = 0, \frac{\partial RSS(\Theta)}{\partial \mu} = 0, \frac{\partial RSS(\Theta)}{\partial \beta_0} = 0, \frac{\partial RSS(\Theta)}{\partial \beta_1} = 0, \frac{\partial RSS(\Theta)}{\partial \beta_2} = 0, \frac{\partial RSS(\Theta)}{\partial \beta_3} = 0, \dots, \frac{\partial RSS(\Theta)}{\partial \beta_p} = 0$. Solution to these equations is difficult since the parameter μ is decenter time point.

We suggest to use the following algorithm while μ is kept fixed.

Algorithm 1

1. Set $\mu = j; j = 1, 2, \dots, N$.
2. For a given μ , obtain LSEs of θ_1, θ_2 using the data $\{(y_1, t_1), (y_2, t_2), \dots, (y_j, t_j)\}$.
3. For a given μ , obtain LSEs of $\beta_0, \beta_1, \beta_2, \beta_3, \dots, \beta_p$ using the data $\{(y_{j+1}, t_{j+1}), (y_{j+2}, t_{j+2}), \dots, (y_N, t_N)\}$.
4. Compute RSS_j using the estimates of $\{\theta_1, \theta_2, \beta_0, \beta_1, \beta_2, \beta_3, \dots, \beta_p\}$ and fixed μ
5. Repeat steps 2-4 for all $j = 1, 2, \dots, N$ and obtain the RSS for each iteration.
6. Search j^* which is a j that corresponds to the minimum RSS.
7. Take $\mu = j^*$ as an optimum value of the parameter μ .

In order to find the optimal value of μ , i.e. the turning point between the exponential and polynomial growth, we will use the technique of minimizing the residual sum squares in “Analysis of COVID-19 Cases in India”.

We will use MAPE (Mean Absolute Percentage Error) to evaluate the performance of the mode.

$$MAPE = \frac{100\%}{N} \sum_{t=1}^N \left| \frac{y_t - \hat{y}_t}{y_t} \right|,$$

where y_t is the observed value at time point t and \hat{y}_t is an estimate of y_t .

In order to make the results easy to interpret, we will also use Accuracy (%).

$$Accuracy (\%) = 100 - MAPE (\%).$$

ARIMA Model

The Auto-Regressive Integrated Moving Averages method gauges the strength of one dependent variable relative to other changing variables. It is one of the most used time-series models in diverse fields of data analysis as it takes into account the changing of trends, periodic changes as well as random disturbances in the time-series data. It is used for both better understanding of the data as well as forecasting, see Brockwell et al. (1996).

Autoregressive model (AR) is effectively merged with the Moving Averages model (MA) to formulate a useful time-series model, ARIMA model. The *Autoregression (AR)* element of the model shows a changing variable that regresses on its own prior values and the *Moving Average (MA)* element incorporates the dependency between an observation and a residual error from a moving average model applied to prior observations. However, this model can only be applied to stationary data. Since many real-life datasets consist of an element of non-stationarity, to model such datasets, ARIMA model was developed. This model is open for non-stationary data as the *Integrated (I)* factor of the model represents the

differencing of raw observations to allow the time-series to become stationary.

Here, we may refer the reader to follow Box et al. (2008, 2015) for more details on ARIMA model, estimation and its application.

The general forms of $AR(p)$ and $MA(q)$ models can be, respectively, represented as the following equations:

$$Y_t = \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \phi_3 Y_{t-3} + \dots + \phi_p Y_{t-p} + \epsilon_t, \tag{3}$$

$$Y_t = \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \theta_3 \epsilon_{t-3} + \dots + \theta_q \epsilon_{t-q} + \epsilon_t, \tag{4}$$

where ϕ s and θ s are auto-regressive and moving averages parameters, respectively, Y_t represents value of time-series at time point t , ϵ_t represents the random disturbance at time point t and is assumed to be independently and identically distributed (i.i.d.) with mean 0 and variance σ^2 .

The ARMA(p, q) model can be represented as:

$$Y_t = \alpha + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \phi_3 Y_{t-3} + \dots + \phi_p Y_{t-p} + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \theta_3 \epsilon_{t-3} + \dots + \theta_q \epsilon_{t-q} + \epsilon_t, \tag{5}$$

where α is an intercept.

The differenced stationary time-series can be modelled as an ARMA model to use ARIMA model on the time-series data (Ceylan 2020; He and Tao 2018; Manikandan et al. 2016). The ARIMA model is generally denoted as ARIMA(p, d, q) where, p is the order of auto-regression, d is the degree of difference and q is the order of moving average.

The degree of difference, i.e. d is a transformation (operator) that is used to make the time-series stationary as it removes the increasing trends. A higher value of d indicates positive autocorrelations out to a high number of lags.

The first step to model the time-series by ARIMA is to determine the time-series data for stationarity. The Augmented Dickey–Fuller (ADF) test may be applied to determine if the time series after differencing is stationary or not.

The ADF test is applied to test the null hypothesis for the presence of a unit root (which indicates non-stationarity of the series).

In order to deduce the $ARIMA(p, d, q)$ model, we can proceed as follows:

We have the $ARMA(p', q)$ represented as follows (as per Eq. 5)

$$Y_t - \phi_1 Y_{t-1} - \phi_2 Y_{t-2} - \phi_3 Y_{t-3} - \dots + \phi_{p'} Y_{t-p'} = \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \theta_3 \epsilon_{t-3} + \dots + \theta_q \epsilon_{t-q} + \epsilon_t$$

It can be equivalently written as:

$$\left(1 - \sum_{i=1}^{p'} \phi_i L^i\right) Y_t = \left(1 + \sum_{i=1}^q \theta_i L^i\right) \epsilon_t, \tag{6}$$

where L is the lag operator, such that- $L^a(Y_t) = Y_{t-a}$.

Now, assume that the polynomial $(1 - \sum_{i=1}^{p'} \phi_i L^i)$ has a unit root (i.e. a factor of $(1 - L)$) of multiplicity d . Then, Eq. (6) can be re-written as:

$$\left(1 - \sum_{i=1}^{p'-d} \phi_i L^i\right) (1 - L)^d Y_t = \left(1 + \sum_{i=1}^q \theta_i L^i\right) \epsilon_t.$$

Or,

$$\left(1 - \sum_{i=1}^{p'} \phi_i L^i\right) (1 - L)^d Y_t = \left(1 + \sum_{i=1}^q \theta_i L^i\right) \epsilon_t, p' = p - d. \tag{7}$$

This can be generalized as:

$$\left(1 - \sum_{i=1}^p \phi_i L^i\right) (1 - L)^d Y_t = \delta + \left(1 + \sum_{i=1}^q \theta_i L^i\right) \epsilon_t, \tag{8}$$

This defines an $ARIMA(p, d, q)$ process with drift $\frac{\delta}{1 - \sum \phi_i}$.

The second step is to plot the graphs of the Autocorrelation function (ACF) and the Partial Autocorrelation Function (PACF) to determine the most-likely values of p and q .

The final step is to obtain the optimal values of p, d and q using the AIC (Akaike Information Criterion), for more details see https://en.wikipedia.org/wiki/Akaike_information_criterion. These information criteria may be used for selecting the best-fitted models. Lower the values of criteria, higher will be its relative quality. The AIC is given by:

$$AIC = -2(\mathcal{L}) + 2K,$$

where K is the number of model parameters, \mathcal{L} is the maximized value of log – likelihood function.

Exponential Smoothing

Exponential smoothing is one of the simple techniques to model time-series data where the past observations are assigned weights that are exponentially decreasing over time. We propose the following models, for modelling of COVID-19 cases [see Holt (1957) and Winters (1960)].

For single exponential smoothing, let the raw observations be denoted by $\{y_t\}$ and $\{s_t\}$ denote the best estimate of trend at time t . Then, $s_0 = y_0, s_t = \alpha y_t + (1 - \alpha)(s_{t-1})$, where $\alpha \in (0,1)$ denotes the data smoothing factor.

For double exponential (Holt–Winters) smoothing, let the raw observations be denoted by $\{y_t\}$, smoothed values $\{s_t\}$, and $\{b_t\}$ denotes the best estimate of trend at time t . Then,

$$s_1 = y_1,$$

$$b_1 = y_2 - y_1,$$

$$s_t = \alpha x_t + (1 - \alpha)(s_{t-1} - b_{t-1}),$$

$$b_t = \beta(s_t - s_{t-1}) + (1 - \beta)b_{t-1},$$

where $\alpha \in (0,1)$ denotes the data smoothing factor and $\beta \in (0,1)$ denotes the trend smoothing factor. For the forecast at $t = (N + m)$ days, (F_{N+m}) is calculated by

$$F_{N+m} = s_t + mb_t.$$

Analysis of COVID-19 Cases in India

For this study, we have used the data available at GitHub, provided by Centre for Systems Science and Engineering (CSSE) at John Hopkins University (see https://github.com/CSSEGISandData/COVID-19/blob/master/csse_covid_19_data/csse_covid_19_time_series/time_series_covid19_confirmed_global.csv). For this study, we use R software. (see R Core Team 2020).

Exponential–Polynomial Regressions

We have used the data from March 4 to July 11 for continuity of the data.

We know that at the beginning of the spread of the disease in India, the growth was exponential and after some time, it was shifted to polynomial. We first obtain optimum turning point of the growth, i.e. when did the growth rate of the disease shifted to polynomial regime from the exponential. We consider both quadratic and cubic regression model for

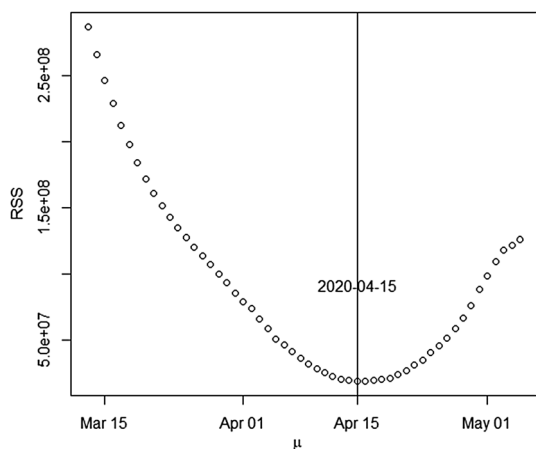


Fig. 1 Trend of RSS and optimum μ for exponential-quadratic regression model

Table 1 Turning point of growth curve for cubic and quadratic regression beyond change point using the COVID-19 cases from 4th March to a given day

Day	Change point	
	Cubic	Quadratic
25th April	5th April	5th April
30th April	5th April	5th April
2nd May	5th April	5th April
3rd May	7th April	5th April
5th May	10th April	11th April
10th May	10th April	18th April

Table 2 Regression table for region I (Exponential Regression)

Parameter	Coefficients	S.E	t	PV
θ_1	16.543	1.969	8.40	$1.7e-09$
θ_2	0.163	0.00389	41.97	$2e-16$

Table 3 Regression models fitting for Region II (5th April–2nd May)

Model	Parameter	OLS Estimates	S.E	t	PV	R^2	(F-statistic, PV)
Linear	β_0	-43,427.02	1889.22	-22.99	$<2e-16$	0.9773	(1119, $<2e-16$)
	β_1	1326.49	39.66	3.45	$<2e-16$		
Quadratic	β_0	17,410.66	1501.022	11.60	$2.52e-11$	0.9997	(3.901e+04, $<2.26e-16$)
	β_1	-1335.45	65.087	-20.52	$<2e-16$		
	β_2	28.32	0.6905	41.01	$<2e-16$		
Cubic	β_0	6196.57	10,073.87	0.615	0.545	0.9997	(2.63e+04, $<2e-16$)
	β_1	-594.89	661.096	-0.9	0.378		
	β_2	12.29	14.2489	0.863	0.397		
	β_3	0.113	0.1009	1.126	0.272		

second part of the data. We will also discuss the types of polynomial growth (with their equations) in India.

In order to find the turning point of the growth curve, we follow the Algorithm 1, given in the previous section. Using that, we evaluate the RSS for all the days (from March 4) and find the date on which it is minimum. The change points of growth curve for cubic and quadratic regressions are presented in Fig. 1 depending upon the size of the data set. From Fig. 1, we can confirm that the growth rate of COVID-19 cases was exponential till April 5 and then after it follows the polynomial growth regime while we use the COVID-19 cases till July 11 (Table 1).

We call the region of exponential growth in India as Region I. The coefficients of the model are presented in Table 2.

We see that after the exponential regime (till April 5), the growth curve follows a polynomial growth till May 2. After this, we again see a change in the behavior of the growth curve. In Tables 3, 4, 5 and 6, we try to model these growth curves through regression analysis.

Having evaluated the coefficients for various models (i.e. linear, quadratic and cubic) as well as the important statistics (i.e. R^2 values, p values of the models as well as individual coefficients and F-statistic), we will select the best-fitting models. In order to select the best-fitting models for Region II (April 6 to May 2); III (May 3 to May 15), IV (May 16 to May 31) and V (June 1 to July 11), we have the following steps. We select that model which has high R^2 values, significant p value, high F-statistic and where the p values of all the variables are significant.

We see for Region II, from Table 3, that the linear model is having a relatively lower F-statistic and R^2 values in comparison to the Quadratic and Cubic models. So, we eliminate the possibility of linear fitting. Further, we see that the p values, F-statistics and the R^2 values are quite significant in both Quadratic as well as the Cubic models. But, if we look at the individual p values of the coefficients, we see that the individual p values are not significant for the Cubic model. On the other hand, the individual p values are significant for

Table 4 Regression models fitting table for Region III (3rd May–15th May)

Model	Parameter	OLS Estimates	S.E	<i>t</i>	PV	<i>R</i> ²	(F-statistic, PV)
Linear	β_0	-2,998,284.88	3244.99	-91.62	$<2e-16$	0.9990	(1.264e+04, $<2e-16$)
	β_1	3584.12	32.11	111.63	$<2e-16$		
Quadratic	β_0	-23,424.15	56,089.18	-0.418	0.68505	0.9997	(1.329e+04, $<2.26e-16$)
	β_1	-1866.15	1111.75	-1.679	0.12416		
	β_2	26.98	5.50	4.903	0.00062		
Cubic	β_0	-8.967e+05	1.837e+06	-0.488	0.637	0.9997	(1.187e+04, $<2e-16$)
	β_1	2.411e+04	5.464e+04	0.441	0.669		
	β_2	-2.305e+02	5.413e+02	-0.426	0.680		
	β_3	8.497e-01	1.786e+00	0.476	0.646		

Table 5 Regression models fitting table for Region IV (16th May–31st May)

Model	Parameter	OLS estimates	S.E	<i>t</i>	PV	<i>R</i> ²	(F-statistic, PV)
Linear	β_0	-403,709.8	9765.5	-41.34	4.92e-16	0.9951	(3069, $<2.2e-16$)
	β_1	6627.8	119.6	55.40	$<2e-16$		
Quadratic	β_0	296,345.02	50,252.76	5.897	5.26e-05	0.9997	(2.285e+04, $<2.26e-16$)
	β_1	-10,606.59	1235.97	-8.582	1.03e-06		
	β_2	105.73	7.58	13.948	3.37e-09		
Cubic	β_0	-7.339+05	1.025e+06	-0.716	0.488	0.9997	(1.525e+04, $<2e-16$)
	β_1	2.746e+04	3.783e+04	0.726	0.482		
	β_2	-3.623e+02	4.649e+02	-0.779	0.451		
	β_3	1.914e+00	1.901e+00	1.007	0.334		

Table 6 Regression models fitting Table for Region V (1st June – 11th July)

Model	Parameter	OLS estimates	S.E	<i>t</i>	PV	<i>R</i> ²	(F-statistic, PV)
Linear	β_0	-1,282,158.3	46,562.7	-27.54	$<2e-16$	0.9725	(1417, $<2.2e-16$)
	β_1	15,845.4	420.9	37.65	$<22e-16$		
Quadratic	β_0	1.669e+06	5.836e+04	28.60	$<22e-16$	0.9996	(4.89e+04, $<2.26e-16$)
	β_1	-3.844e+04	1.070e+03	-35.94	$<22e-16$		
	β_2	2.468e+02	74.856e+00	50.81	$<22e-16$		
Cubic	β_0	1.786+06	2.465e+05	-7.245	1.342e-08	0.9999	(2.019e+05, $<2e-16$)
	β_1	5.711e+04	6.799e+03	8.400	4.242e-10		
	β_2	-6.280e+02	6.215e+01	-10.104	3.462e-12		
	β_3	2.651e+00	1.8822e-01	14.082	$<22e-16$		
Quartic	β_0	-1.168e+07	2.067e+06	-5.649	2.082e-06	1.000	(2.419e+05, $<2e-16$)
	β_1	4.223e+05	7.617e+04	5.545	2.822e-06		
	β_2	-5.658e+03	1.048e+03	-5.401	4.392e-06		
	β_3	3.329e+01	6.375e+00	5.221	7.632e-06		
	β_4	-6.963e-02	1.4492e-02	-4.807	2.712e-05		
Quintic	β_0	-5.784e+07	2.094e+07	-2.762	0.00908	1.000	(2.145e+05, $<2.2e-16$)
	β_1	2.554e+06	9.651e+05	2.646	0.01213		
	β_2	-4.486e+04	1.773e+04	-2.530	0.01606		
	β_3	3.923e+02	1.623e+02	2.418	0.02095		
	β_4	-1.709e+00	7.4022e-01	-2.308	0.02702		
	β_5	2.980e-03	1.3462e-03	2.214	0.03341		

Table 7 Parameters for exponential regression in Region V

Parameter	OLS Estimate	S.E	<i>t</i>	<i>p</i> value
α	8.718e+03	1.385e+02	62.95	$< 2e-16$
β	3.530e-02	1.333e-04	264.69	$< 2e-16$

the Quadratic model. Thus, we can conclude that the *Quadratic* model is the best-fitting model for Region II (April 6 to May 2).

For Region III, from Table 4, that all the three models have high F-statistic values, high *p* values and high R^2 values. But we notice that the coefficient individual *p* values are not significant in both Quadratic and Cubic models. Thus, we conclude that the *Linear* model is the best-fitting model for Region III (May 3 to May 15).

For Region IV, from Table 5, we see that the R^2 values for all the models are very high. All the models also have significant *p* values. The F-statistic of both Quadratic and Cubic models are also high. But, the coefficient individual *p* values are not significant in the Cubic model. Thus, we conclude that the *Quadratic* model is the best-fitting model for Region IV (May 16 to May 31).

For Region V, from Table 6, we see that the R^2 values of all the models are very high (quadratic, cubic, quartic and quintic models have exceptionally high). All the models also have significant *p*-values. The F-statistic of Quadratic, Cubic, Quartic and Quintic models is high. F-statistic value of Quartic model is the highest. The coefficient individual *p* values of Quartic model are also significant. Thus, we conclude that the *Quartic* model is the best-fitting model for Region V (June 1 to July 11).

Note For Region V, due to spike in the cases, we also checked the fitting of exponential curve in this region (Table 7).

Let the exponential model be- $y(t) = \alpha e^{\beta t}$

We obtained the following parameter values:

The RSE for this model is 4178 and MAPE (%) is 0.85%. Both of these values are quite larger than those of Quartic model (Refer Table 9 for RSE and MAPE values of Quartic model in Region V). Thus, we conclude that Quartic model is the best-fitting model for Region V (1st June to 11th June).

All the ANOVA tables (Refer to Table 8) for Region II, III, IV and V suggest significant *p*-values for its coefficients and suggest that the models fit well the respective regions.

Thus, according to our study, the growth of the virus was exponentially increasing from March 4 to April 5. Then

Table 8 ANOVA table for Region II (Quadratic Regression), III (Linear Regression), IV (Quadratic Regression) and V (Quartic Regression)

Region	Model	Variable	Degrees of freedom	Sum of squares	Mean sum of squares	F-statistic	<i>p</i> value
II	Quadratic	<i>t</i>	1	2,882,180,732	2,882,180,732	76,347	$< 2.2 \times 10^{-16}$
		t^2	1	63,489,615	63,489,615	1681	$< 2.2 \times 10^{-16}$
		Residuals	24	906,026	37,751		
III	Linear	<i>t</i>	1	2,337,950,722	2,337,950,722	12,462	$< 2.2 \times 10^{-16}$
		Residuals	11	2,063,722	187,611		
IV	Quadratic	<i>t</i>	1	1.4935e+10	1.4935e+10	45,505.60	$< 2.2 \times 10^{-16}$
		t^2	1	6.3857e+07	6.3857e+07	194.56	3.373e-09
		Residuals	13	4.2667e+06	3.2821e+05		
V	Quartic	<i>t</i>	1	1.4412e+12	1.4412e+12	941,897.436	$< 2.2 \times 10^{-16}$
		t^2	1	3.9077e+10	3.9077e+10	25,539.326	$< 2.2 \times 10^{-16}$
		t^3	1	4.8467e+08	4.8467e+08	316.761	$< 2.2 \times 10^{-16}$
		t^4	1	3.5353e+07	3.5353e+07	23.105	$< 2.2 \times 10^{-16}$
		Residuals	36	5.5083e+07	1.5301e+06		

Table 9 Course of COVID-19 growth in India (March 4 to July 11)

Region	Dates	Best-fitted model	MAPE (%)	RSE
I	March 4th to April 5th	$y(t) = 16.54 \times e^{0.163t}$	8.60	81.66
II	April 6th to May 2nd	$y(t) = 17410.67 - 1335.45t + 28.32t^2$	0.83	194.3
III	May 3rd to May 15th	$y(t) = -29825 + 3584t$	0.49	433.1
IV	May 16th to May 31st	$y(t) = 296345.02 - 10606.59t + 105.73t^2$	0.31	572.6
V	June 1st to July 11th	$y(t) = -1.168 \times 10^7 + 4.223 \times 10^5t - 5.658 \times 10^3t^2 + 33.29t^3 - 0.06963t^4$	0.21	1237

after, the virus grew by following a quadratic rate from April 6 to May 2. After May 3, we experienced a linear growth. But after May 15 to May 31, we experienced a sudden rise in the rate of growth of the virus and have seen quadratic growth again. Further, for the period of June 1 to July 11, we see experienced a quartic (4-degree polynomial) growth, which is very alarming (see Table 9 for best-fitted regression models). Figure 2 shows the best-fitted regression models

to the daily cumulative cases of COVID-19 in India from March 4 to July 11 (Table 10).

Time-series Models Fitting

We use the daily time-series data of number of cumulative confirmed cases from March 4 to July 10.

First, we check the stationarity of the transformed time-series using ADF Tests. Dickey–Fuller statistic is 6.3915 with p value 0.99 which indicates that the growth of COVID-19 cases is not stationary. The ARIMA models may be useful over the ARMA models. The ACF and PACF plots are shown in Fig. 3.

We then obtain the optimal ARIMA parameters (p, d, q) using the AIC. We take various possible combinations of (p, d, q) and compute the AIC. Then, select the best-fitted ARIMA model that has the lowest AIC among all considered models. According to the AIC, the ARIMA (5, 2, 5) is the best-fitted model for the COVID-19 cases, India (see Table 11). Estimates of ARIMA (5, 2, 5) parameters and MAPE are shown in Table 11.

Table 10 AIC for ARIMA models for COVID-19 cases, India (4th March to 10th July)

Model			AIC
p	d	q	
1	1	0	2064.921
2	1	0	2065.752
3	1	0	2067.442
2	1	2	2063.402
3	1	5	2055.702
2	0	1	2070.555
1	2	1	2043.075
2	2	1	2044.97
1	2	2	2045
2	2	2	2022.933
3	2	2	2048.967
3	2	3	2021.517
5	2	5	2001.998
2	2	5	2012.213
5	2	4	2008.532
3	2	5	2013.83

Interpretation of the Parameters

We have selected the model parameters using the Akaike Information Criterion. We obtained the parameters as: $p = 5, d = 2$ and $q = 5$. As $p = 5$, it means that the order (number of time lags) of Autoregression part of the model

Fig. 2 Fitted regression models to the daily cumulative cases of COVID-19 in India till July 11

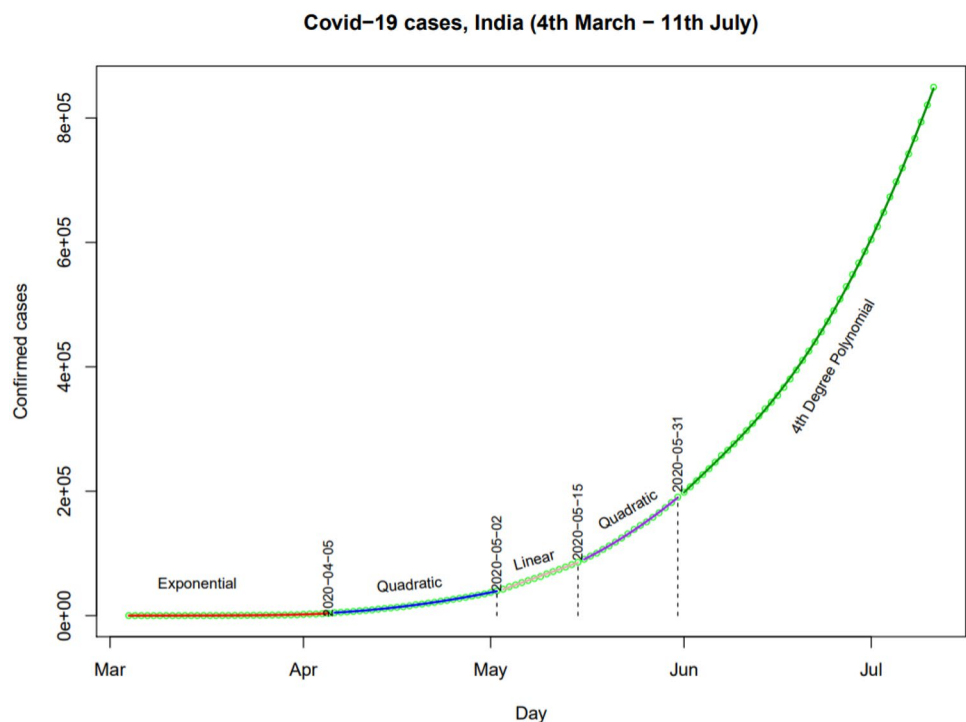


Fig. 3 ACF and PACF for COVID-19 cases in India (4th March to 10th July)

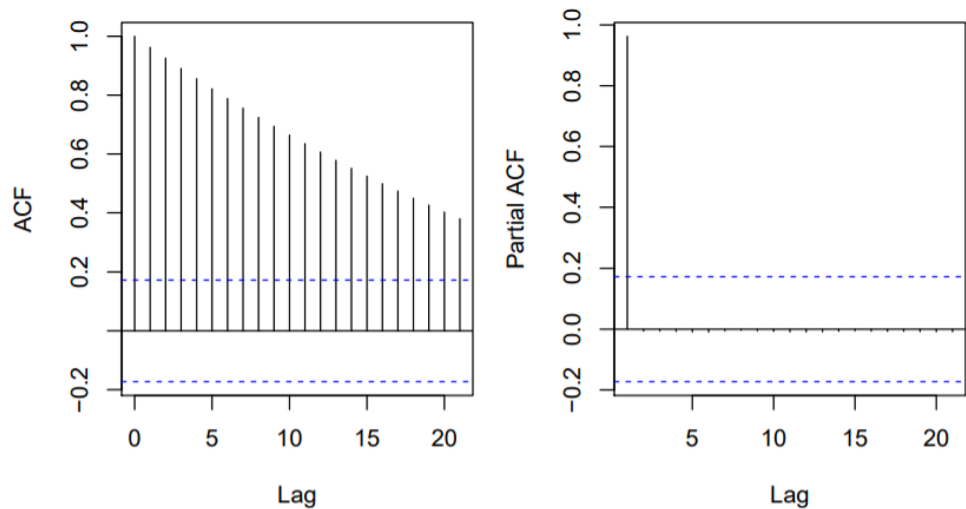


Table 11 Estimates of ARIMA (5,2,5) parameters and MAPE (4th March to 10th July)

Coefficients	Estimate	S.E	MAPE	Accuracy
AR 1	0.6734	0.0566	2.6511%	97.3849%
AR 2	0.3395	0.0773		
AR 3	-0.1541	0.0879		
AR 4	-0.7314	0.0673		
AR 5	0.8639	0.0549		
MA 1	-0.8661	0.0748		
MA 2	-0.4824	0.1081		
MA 3	0.3273	0.1116		
MA 4	0.9873	0.0844		
MA 5	-0.8447	0.0622		

is 5. In general, we can say that the cumulative cases of COVID-19 in a day are dependent on the cases of previous 5 days. As $q = 5$, the present value is dependent on the moving average (residuals) of previous 5 days. As $d = 2$, the series $y_t^* = y_t - 2y_{t-1} + y_{t-2}$ is stationary. A higher value of d indicates positive autocorrelations out to a high number of lags. Thus, we can have the equation for our model, using Eq. 8 as:

$$\left(1 - \sum_{i=1}^5 \varphi_i L^i\right) (1 - L)^2 Y_t = \delta + \left(1 + \sum_{i=1}^5 \theta_i L^i\right) \varepsilon_t,$$

where all the symbols have their meanings as per “ARIMA Model”.

Estimates of the Holt–Winters exponential smoothing and exponential smoothing models are given in Table 12. According to the MAPE and accuracy measures, the ARIMA (5, 2, 5) is a better model than the Holt–Winters exponential smoothing and usual exponential smoothing models. From this, we can conclude that the ARIMA model is the best

Table 12 Estimates and MAPE of exponential smoothing models

Model	Parameter	Estimate	MAPE	Accuracy
Holt–Winters exponential smoothing	α	1	2.8911%	97.1089%
	β	1		
	a	820,916		
	b	271,145		
Exponential smoothing	α	0.9999	7.5384%	92.4616%
	a	820,914.7000		

fit for the cases of COVID-19, followed by Holt–Winters model. The forecasting values along with 95% confidence intervals are shown in Table 13 and Fig. 4. We have used actual data from 11th June to validate the model.

Even though most of the actual cases are covered in the 95% confidence intervals of the ARIMA and Holt–Winters forecasts, they are seen to be nearer to the Upper Limits of the Confidence Intervals and are deviated from the estimates. It might be possible that in the future days, the forecasts might underestimate the actual cases. This might be attributed to the changing pattern of the growth of the pandemic in our country as seen in the regression analysis. Thus, we suggest a segment-wise time-series models to forecast the future cases in a more accurate manner.

We present the segment-wise ARIMA and Holt–Winters models for 1st June to 10th July.

We have seen that our time-series data are non-stationary and, thus, we select the most optimal values of (p, d, q) , which has the least AIC. According to AIC, (5, 2, 3) is the best-fitting model for the time-series data from June 1 to July 10, with AIC = 634.18. Estimates of ARIMA (5, 2, 3) model with the corresponding MAPE and Accuracy are given in Table 14 (Fig. 5).

Table 13 Forecast using ARIMA and holt-winters models for 10 days

Day	ARIMA			Holt-winters			Actual
	Estimate	Lower	Upper	Estimate	Lower	Upper	
11th July	848,374.1	847,247.1	849,501.0	848,030	846,636.1	849,423.9	849,522
12th July	875,792.9	873,473.0	878,112.8	875,144	872,027.2	878,260.8	878,254
13th July	903,010.8	899,507.1	906,514.5	902,258	897,042.5	907,473.5	906,752
14th July	930,330.7	925,622.6	935,038.7	929,372	921,737.3	937,006.7	936,181
15th July	958,494.9	952,436.5	964,553.2	956,486	946,148.6	966,823.4	968,857
16th July	987,619.1	979,891.2	995,347.0	983,600	970,303.1	996,896.9	1,003,832
17th July	1,017,773.7	1,007,937.9	1,027,609.6	1,010,714	994,221.2	1,027,206.8	1,039,084
18th July	1,048,569.9	1,036,227.8	1,060,912.1	1,037,828	1,017,919.2	1,057,736.8	1,077,781
19th July	1,079,470.6	1,064,335.7	1,094,605.5	1,064,942	1,041,410.4	1,088,473.6	1,118,206
20th July	1,110,527.9	1,092,443.6	1,128,612.2	1,092,056	1,064,705.9	1,119,406.1	1,155,338

Fig. 4 Fitted ARIMA (5, 2, 3) and exponential smoothing models and forecasting from ARIMA for Covid-19 cases in India (stars show the actual observations). Model built on data from June 1 to July 10

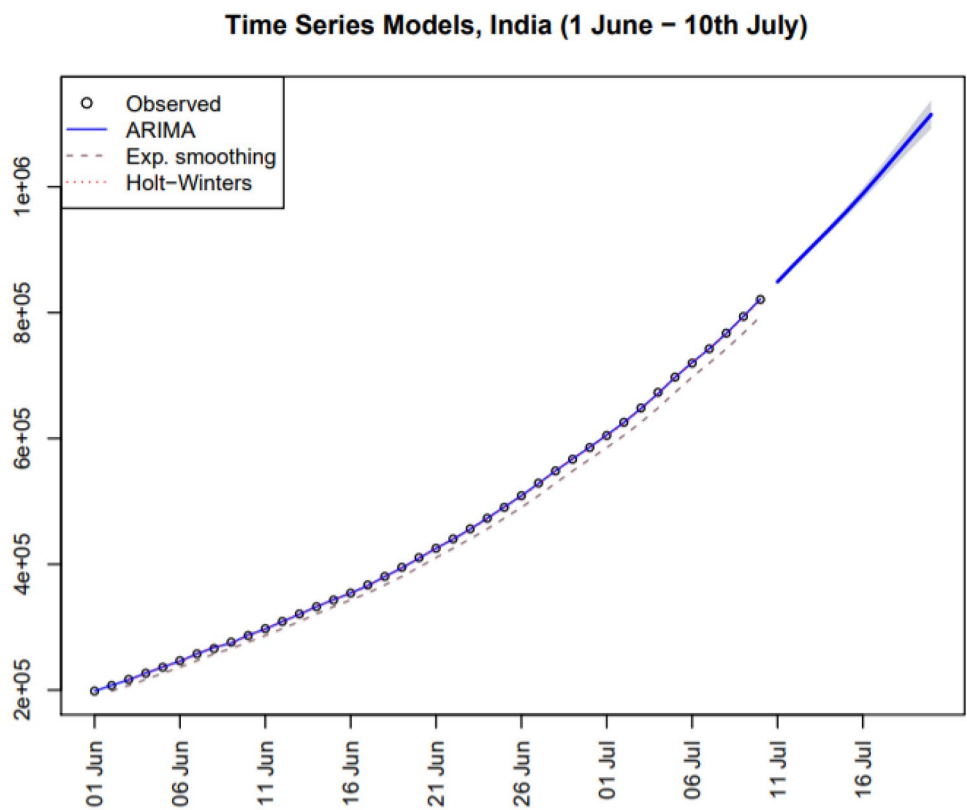


Table 14 Estimates of ARIMA (5,2,3) parameters and MAPE (1 June to 10 July)

Coefficients	Estimate	S.E	MAPE	Accuracy
AR 1	1.9985	0.2455	0.1364%	99.8636%
AR 2	-2.0142	0.4505		
AR 3	1.2380	0.4533		
AR 4	-0.6108	0.3606		
AR 5	0.3742	0.1895		
MA 1	-2.2358	0.3847		
MA 2	2.0805	0.7260		
MA 3	-0.7309	0.3853		

From Fig. 6, we deduce that the optimal value of d is 2, as the time series becomes stationary with differencing degree = 2.

Estimates of the Holt–Winters exponential smoothing and exponential smoothing models are given in Table 15. According to the MAPE and accuracy measures, the ARIMA (5, 2, 3) is a better model than the Holt–Winters exponential smoothing and usual exponential smoothing models. From this, we can conclude that the ARIMA model is the best fit for the cases of COVID-19, followed by Holt–Winters model. The forecasting values along with 95% confidence intervals are shown in Table 16 and Fig. 4. We have used

Fig. 5 ACF and PACF plots for COVID-19 cases (1 June to 10 July)

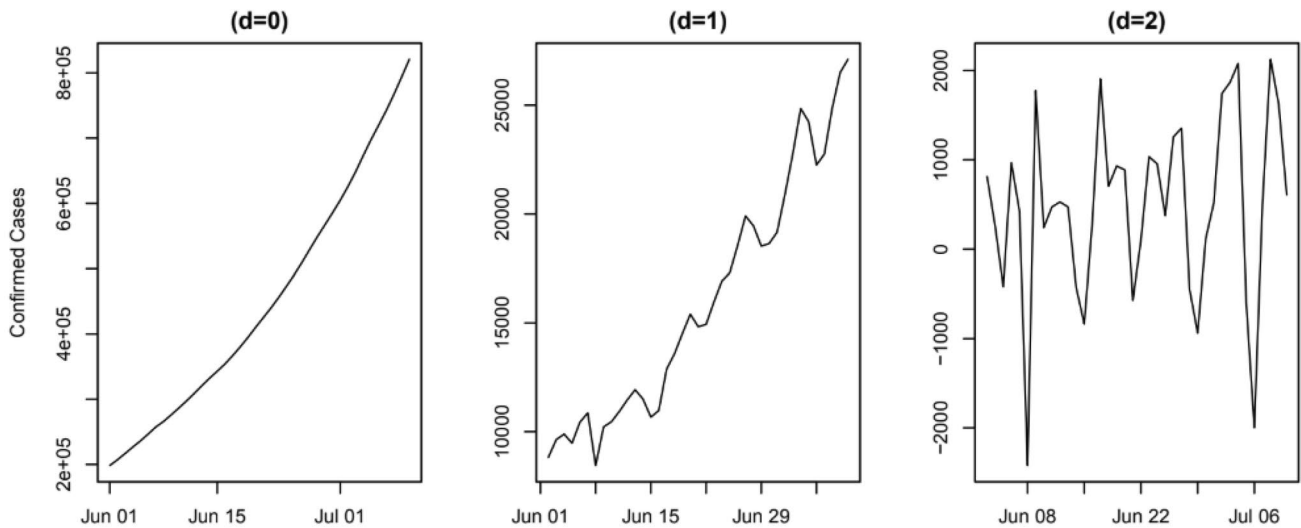
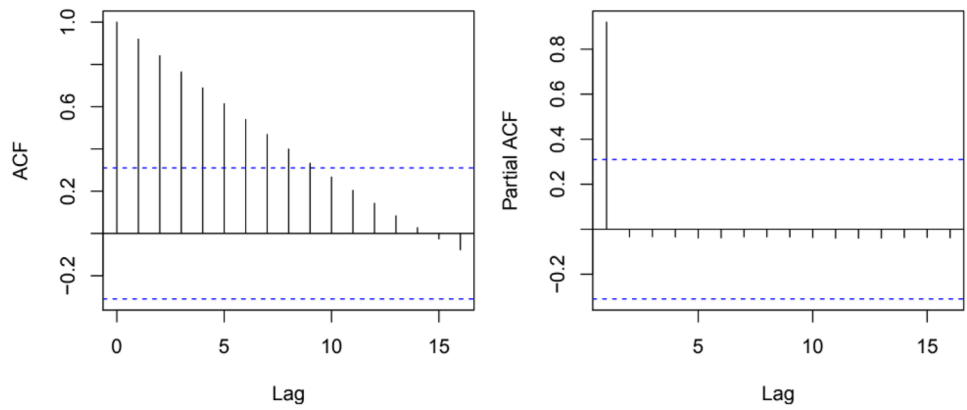


Fig. 6 Difference plots for COVID-19 cases (1 June to 10 July)

Table 15 Estimates and MAPE of exponential smoothing models

Model	Parameter	Estimate	MAPE	Accuracy
Holt–Winters exponential smoothing	α	1	0.2172%	99.7828%
	β	1		
	a	820,916		
	b	271,145		
Exponential smoothing	α	0.9999	3.5757%	96.4243%
	a	820,914.7000		

actual data from 11 June to validate the model. We observe that the ARIMA model captures the trend well and estimates the cumulative cases properly.

From the interpretations of both the fitted ARIMA models, we can say that as the values of p and d are 5 and 2, respectively; the daily cumulative cases are dependent on the

cases of previous 5 days. Also, to convert the time series of daily cases into stationary, we need differencing degree of 2.

Conclusion and Future Scope

From the regression analysis, we conclude that the spread of COVID-19 disease grew exponentially from March 3 to April 5. Further, from April 6 to May 2, the cases followed a quadratic regression. From May 3 to May 15, we see a linear growth of the pandemic with average daily cases of 3584. After May 15 to May 31, we again saw a spike in the cases that lead to a quadratic growth of the pandemic. And, from June 1 to July 11, we saw a major spike in the growth of the pandemic as it has followed quartic growth.

Verma et al. (2020) showed the four stages of the epidemic, S1: exponential, S2: power law, S3: linear and S4: flat. We saw that the course of COVID-19 in India followed this regime till May 15. But after the linear trend from May

Table 16 Forecast using ARIMA and Holt–Winters models for 10 days (Model based on data from June 1 to July 11)

Day	ARIMA			Holt–Winters			Actual
	Estimate	Lower	Upper	Estimate	Lower	Upper	
11th July	849,004.1	847,380.2	850,628.0	848,030	845,992.4	850,067.6	849,522
12th July	877,131.0	873,820.0	880,441.9	875,144	870,587.8	879,700.2	878,254
13th July	904,470.9	899,668.1	909,273.7	902,258	894,634.0	909,882.0	906,752
14th July	931,603.7	925,143.6	938,063.8	929,372	918,211.5	940,532.5	936,181
15th July	959,587.9	951,065.0	968,110.8	956,486	941,374.7	971,597.3	968,857
16th July	989,057.8	977,968.3	1,000,147.4	983,600	964,162.4	1,003,037.6	1,003,832
17th July	1,020,020.5	1,005,736.9	1,034,304.1	1,010,714	986,604.7	1,034,823.3	1,039,084
18th July	1,051,860.3	1,033,742.1	1,069,978.4	1,037,828	1,008,725.1	1,066,930.8	1,077,781
19th July	1,083,688.0	1,061,252.4	1,106,123.6	1,064,942	1,030,543.2	1,099,340.8	1,118,206
20th July	1,114,984.3	1,087,911.9	1,140,256.8	1,092,056	1,052,075.2	1,132,036.8	1,155,338

3 to May 15, the spread has again reached the quadratic growth and from June 1 to July 11, India is witnessing a quartic growth. This might be attributed to the relaxation of lockdown measures in the country. Though it was much likely that the cases would start to reduce post-linear stage growth as the total cases may start to follow a square root equation, i.e. $y(t) \sim \sqrt{t}$. And this might lead to reduction in the daily number of cases ($\text{asy}'(t) \sim 1/\sqrt{t}$), leading to flattening of the curve. But, due to reduced restrictions, we see a reverse trend, which might be alarming and suggest the imposition of strict lockdown to reverse this trend of pandemic growth. If we continue to open our economy in this way, we might go back to the exponential growth of the pandemic and this would lead to huge destruction to human lives and cause a greater impact on our economy.

We also observe that some cities have been the hotspots of the disease, such as Delhi (more than 131,000 cases), Mumbai (more than 94,000 cases), Chennai (more than 78,000 cases), Thane (more than 63,000 cases), etc. as on 14 June, 2020. While the other states and cities have seen a slower growth of the pandemic, these cities have seen explosive growths. Due to the opening of air and rail transport in the country, the virus is likely to spread in the other regions as well as people from these cities (especially metro cities) are travelling to different states. Thus, it is highly advisable that the country should go back to its lockdown phase until we see reduction in trend.

In time-series analysis, we conclude that the ARIMA (5, 2, 5) is the best-fitting model for the cases of COVID-19 from 4th March to 10th July with an accuracy of 97.38%. The basic exponential smoothing is not very accurate for our case, but we see that the Holt–Winters model is around 97.11% accurate. Both ARIMA (5, 2, 5) and Holt–Winters models suggest a rise in the number of cases in the coming days. We observed that both the ARIMA and Holt–Winters models capture the data well and the actual data from 11th July validate the forecasts well as they lie in the predicted confidence intervals. But, while validating the model, the

actual values are always near to the Upper Confidence Limits, it might be possible that in further days, our model might underestimate the cases. This might be possible because of the changing trend of the growth of the pandemic in India.

Thus, we used segmented time-series models and took data from 1st June to 10th July to build separate ARIMA and Holt–Winters models. We concluded that ARIMA (5, 2, 3) is the best-fitting model for COVID-19 cases in the given time period with an accuracy of 99.86%. The basic exponential smoothing is not very accurate or this case as well but, the Holt–Winters model is around 99.78% accurate. We also observe that the ARIMA and Holt–Winters models capture the data well and the actual data from 11th July validate the forecasts and lie near to the estimates.

We may also conclude that the cases of COVID-19 will rise in the coming days and the situation may turn alarming if proper measures are not followed. Since the economic activities have started in the country, people need to be more careful while going out. And explosion of the pandemic in the whole country can cause a serious damage to human lives, healthcare system as well as the economy of the country. Thus, there is an urgent need of imposing strict lockdown measures to curb the growth of the pandemic. We must also learn to lead our lives by following all the precautions even if the lockdown restrictions are relaxed and the economic activities are resumed.

Comparison of Indian scenario with that of other countries might not prove fruitful at this stage because of the demographic differences and/or the characteristics of the disease. Also, comparison of the Indian context with that of the other countries of the world will require to study the spread of the pandemic in those countries in depth and might be considered as an altogether in the future studies.

This study was limited to data-driven models using the total COVID-19 cases. In the future studies, the other co-factors (associated with the demographics, social, cultural and medical infrastructure, etc.) can be taken to considerations.

Acknowledgements Dr. Vikas Kumar Sharma greatly acknowledges the financial support from Science and Engineering Research Board, Department of Science & Technology, Govt. of India, under the scheme Early Career Research Award (file no.: ECR/2017/002416).

References

- Arti MK, Kushagra B (2020) Modeling and predictions for COVID 19 spread in India, <https://doi.org/10.13140/RG.2.2.11427.81444>
- Bhatnagar MR (2020) COVID-19: mathematical modeling and predictions, submitted to ARXIV. Online available at: <https://web.iitd.ac.in/~manav/COVID.pdf>
- Box GEP, Jenkins GM, Reinsel GC (2008) Time analysis. Wiley, Hoboken
- Box GEP, Jenkins GM, Reinsel GC, Ljung GM (2015) Time series analysis: forecasting and control. Wiley, Hoboken
- Brockwell PJ, Davis RA (1996) Introduction to time series and forecasting. Springer, Berlin
- Ceylan Z (2020) Estimation of COVID-19 prevalence in Italy, Spain, and France. *Total Environ Sci Total Environ* 729:138817
- Chatterjee S, Shayak B, Asad A, Bhattacharya S, Alam S, Verma MK (2020) Evolution of COVID-19 pandemic: power law growth and saturation. *medRxiv*. <https://doi.org/10.1101/2020.05.05.20091389>
- Chintalapudi N, Gopi B (2020) Amenta Francesco COVID-19 virus outbreak forecasting of registered and recovered cases after sixty day lockdown in Italy: a data driven model approach. *J Microb Immunol Infect*. <https://doi.org/10.1016/j.jmii.2020.04.004>
- Earnest A, Chen MI, Ng D, Leo YS (2005) Using autoregressive integrated moving average (ARIMA) models to predict and monitor the number of beds occupied during a SARS outbreak in a tertiary hospital in Singapore. *BMC Health Serv Res* 5:1–8. <https://doi.org/10.1186/1472-6963-5-36>
- Fanelli D, Francesco P (2020) Analysis and forecast of COVID-19 spreading in China, Italy and France. *Chaos Solitons Fractals Nonlinear Sci Nonequilibrium Complex Phenomena* 134:109761
- Gaudart J, Touré O, Dessay N, Dicko AL, Ranque S, Forest L, Demongeot J, Doumbo OK (2009) Modelling malaria incidence with environmental dependency in a locality of Sudanese savannah area. *Mali Malar J*. <https://doi.org/10.1186/1475-2875-8-61>
- Giordano G, Blanchini F, Bruno R, Colaneri P, Filippo A, Matteo A (2020) Force, a SIDARTHE model of COVID-19 epidemic in Italy. *arXiv:2003.09861* [q-bio.PE]
- Gupta PK, Bhaskar P, Maheshwari S (2020) Coronavirus 2019 (COVID-19) Outbreak in India: a perspective so far. *J Clin Exp Invest*. 11(4):em00744
- He Z, Tao H (2018) International journal of infectious diseases epidemiology and ARIMA model of positive-rate of influenza viruses among children in Wuhan, China: a nine-year retrospective study. *Int J Infect Dis* 74:61–70. <https://doi.org/10.1016/j.ijid.2018.07.003>
- Holt CE (1957) Forecasting seasonal and trends by exponentially weighted averages (O.N.R. Memorandum No. 52). Carnegie Institute of Technology, Pittsburgh. <https://doi.org/10.1016/j.ijfor-ecast.2003.09.015>
- Ivorra B, Ferrández MR, Vela-Pérez M, Ramos AM (2020) Mathematical modeling of the spread of the coronavirus disease 2019 (COVID-19) considering its particular characteristics. The case of China. *Commun Nonlinear Sci Numerical Simulat* 88:105303
- Kumar A, Gupta PK, Srivastava A (2020) A review of modern technologies for tackling COVID-19 pandemic. *Diabetes Metab Syndr* 14(4):569–573
- Lin Q, Zhao S, Gao D, Lou Y, Yang S, Musa SS, Wang MH, Cai Y, Wang W, Yang L, He D (2020) A conceptual model for the coronavirus disease 2019 (COVID-19) outbreak in Wuhan, China with individual reaction and governmental action. *Int J Infect Dis* 93:211–216
- Manikandan M, Velavan A, Singh Z, Purty AJ, Bazroy J, Kannan S (2016) Forecasting the trend in cases of Ebola virus disease in West African countries using auto regressive integrated moving average models. *Int J Community Med Public Health* 3:615–618
- Ministry of Health and Family Welfare (2020) Government of India. Available at <https://www.mohfw.gov.in/>
- Montgomery DC, Peck EA, Vining GG (2012) Introduction to linear regression analysis. Wiley, Hoboken
- Petropoulos F, Makridakis S (2020) Forecasting the novel coronavirus COVID-19. *PLoS ONE* 15(3):e0231236. <https://doi.org/10.1371/journal.pone.0231236>
- Polwiang S (2020) The time series seasonal patterns of dengue fever and associated weather variables in Bangkok (2003–2017). *BMC Infect Dis* 20:208. <https://doi.org/10.1186/s12879-020-4902-6>
- R Core Team (2020) R: language and environment for statistical computing. R foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>
- Singh B, Jadaun GS (2020) Modeling tempo of COVID-19 pandemic in India and significance of lockdown. *Medrxiv*. <https://doi.org/10.1101/2020.05.15.20103325>
- Verma MK, Asad A, Chatterjee S (2020) COVID-19 pandemic: power law spread and flattening of the curve. *Trans Indian Natl Acad Eng*. <https://doi.org/10.1007/s41403-020-00104-y>
- Winters PR (1960) Forecasting sales by exponentially weighted moving averages. *Manage Sci* 6:324–342
- Zhang X, Liu Y, Yang M, Zhang T, Young AA, Li X (2013) Comparative study of four time series methods in forecasting typhoid fever incidence in China. *PLoS ONE*. <https://doi.org/10.1371/journal.pone.0063116>
- Zhang X, Ma R, Wang L (2020) Predicting turning point, duration and attack rate of COVID-19 outbreaks in major Western countries. *Chaos Solitons Fractals Nonlinear Sci Nonequilibrium Complex Phenomena* 135:109829
- Ziff AL, Ziff RM (2020) Fractal kinetics of COVID-19 pandemic. *medRxiv*. <https://doi.org/10.1101/2020.02.16.20023820>

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.