



Data Smoothing and Other Methods for Generating Forecasts for COVID-19 Cases in Sri Lanka

*¹Gethmini Siriwardena, ²Gayan Dharmaratne, ³Dhammika Amaratunga

^{1,2}Department of Statistics, Faculty of Science, University of Colombo, Sri Lanka

³Princeton Data Analytics, NJ, USA

Email address of the corresponding author - * getsiriwardena@gmail.com

ARTICLE INFO

Article History:

Received: 10 September 2023

Accepted: 01 November 2023

Keywords:

Arima; Smoothing; Trend analysis; Time series

Citation:

Gethmini Siriwardena, Gayan Dharmaratne, Dhammika Amaratunga. (2023). Data Smoothing and Other Methods for Generating Forecasts for COVID-19 Cases in Sri Lanka. Proceedings of SLIIT International Conference on Advancements in Sciences and Humanities, 1-2 December, Colombo, pages 253-258.

ABSTRACT

The COVID-19 pandemic has significantly impacted global society, including Sri Lanka, necessitating the need for reliable forecasting methods. This study compares ten distinct models to predict the number of confirmed COVID-19 cases in Sri Lanka, aiming to assess the performance of statistical models using limited and volatile real-world data characterized by trends, random peaks, and autocorrelations. In addition to the classical ARIMA model, various smoothing and filtering techniques were explored to capture the unique characteristics of the data. The model consistencies in multiple-day predictions were demonstrated, and robust evaluation criteria, along with non-robust measures, were utilized to enhance the effectiveness of the evaluation process. The results highlight the effectiveness of traditional smoothing strategies such as Simple Exponential Smoothing, Holt's Exponential Smoothing, and the Smoothing Splines technique coupled with the ARIMA model. Notably, applying the ARIMA model directly to the original data without smoothing or filtering approaches yielded inadequate forecasts, underscoring its limitations in volatile data settings.

1. INTRODUCTION

The COVID-19 pandemic has had a profound and lasting impact on daily life worldwide, including Sri Lanka. Given the need to implement preventive measures to mitigate the impact of the pandemic, a reliable forecasting approach becomes crucial. It helps public health experts allocate resources, prepare for outbreaks, establish healthcare systems, and make decisions on lockdowns, and other mitigation plans. However, available data, such as the data set that is used in the current study is a relatively short time series which are essentially dominated by significant autocorrelations, trends, and outliers, rather than aspects like seasonal variations. Therefore, standard time series methodologies have some challenges in performing well. Therefore, this study aims to compare the forecasting performance of alternative statistical models for predicting COVID-19 cases in Sri Lanka, offering insights into optimal methods for diverse circumstances.

In general, COVID-19 data are characterized by dominant trends. Smoothing and filtering techniques can effectively address these issues in time series data. For trend capture, LOWESS (Cleveland, 1979) and spline (Craven and Wahba, 1979) models have been widely used. LOWESS smoothing has been applied to represent the trend in a study of tetanus vaccines and TBE, (Hainz et al., 2005), while cubic spline smoothing based on a stochastic state space model has shown promising results in predicting COVID-19 cases and fatalities (Gecili et al., 2021). Spline smoothing techniques have also been used in conjunction with machine learning methods for effective forecasting of COVID-19 cases (Amaratunga, et al, 2023). Exponential smoothing methods, such as Holt's model, have also been commonly employed and found to be reliable for short-term forecasts of daily COVID-19 cases (Martinez et al., 2020). Furthermore, time series data characterized by zero-inflation (Tawiah et al., 2021), was analysed

in modelling of COVID-19 deaths in Ghana.

However, there is a need for further research on combining ARIMA with other filters or smoothers for improved predictions, as well as a scarcity of literature specifically focusing on COVID-19 cases and fatalities in Sri Lanka.

2. MATERIALS AND METHODS

The dataset used in this study was sourced from the website <https://ourworldindata.org/>. The variable, namely 'new cases', was utilized to compare different models for predicting COVID-19 cases in the context of Sri Lanka. The data collection period began on January 27, 2020, and was updated daily. The study covers a time span of over two years until the decline of COVID-19 cases in Sri Lanka, comprising a total of 892 data points until July 6, 2022. The time series considered in this study is illustrated in the accompanying figures.

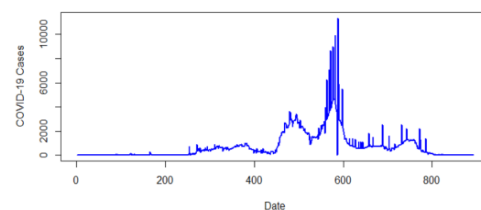


Figure 1 - Plot of COVID-19 Cases

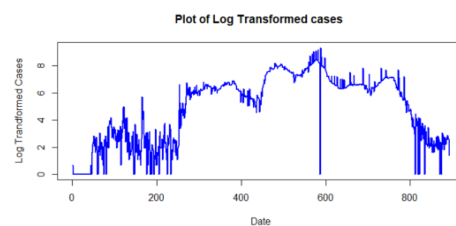


Figure 2 - Plot of Log Transformed COVID-19 Cases

The COVID-19 trajectory for cases exhibits an increasing trend from low levels, followed by a decline towards zero after reaching a peak. Moreover, there are no strong seasonal patterns observed in the plot or in the ACF analysis.

However, fragmentary trends and significant autocorrelations at higher lags suggest the presence of short-term correlations and non-stationarity. Notably, there are prominent outliers, emphasizing the importance of robust evaluation metrics. Considering these characteristics, trend components and sporadic peaks dominate the series, potentially posing challenges for statistical modelling. Nevertheless, the application of smoothing and filtering techniques may help address these complexities. The time series plot exhibits significant variability, outliers, and extended tails. To mitigate these issues, a log transformation was applied, effectively reducing skewness and variability in the data. However, before comparing between models, all the models were converted back to the original scale. Since there were days with no cases (65 days), when the log transformation was applied, an extra 1 was added to all the cases in each day.

To achieve the objective of comparing methods that can capture the properties of the data and to evaluate their forecasts, in this study, the training data set is continuously updated with each realization of the testing data, ensuring adaptiveness and continuity throughout the analysis. For example,

	Range of fitting	Predicting	
Next day predictions	1:x	x+1	x;600,601,602,...,891
One-month (4-weeks) ahead predictions	1:x	x+28	x;600,601,602,...,864

2.1. MODEL IMPLEMENTATION

Several techniques were utilized in this study for comparison purposes:

01. ARIMA on the original data (RAW)- First, the ARIMA model was fitted on the raw data as a baseline model to gauge how different (or

better) the other models performed against it. The `auto.arima()` function in R was used to fit the ARIMA model. ARIMA (3,1,2) was fitted on raw cases.

02. ARIMA on the log transformed data (LOG) - The log transformation was applied to all other models in this analysis, excluding the previous case. The ARIMA model was also fitted to this, with ARIMA (1,1,1) being the best fit for log-transformed cases.

03. ARIMA on the log transformed, LOWESS (linear fit) smoothed data. (LOWESS)- The LOWESS smoother was applied to the logged data set using the `'lowess()'` function in R. A two-week moving window was selected as the smoother's bandwidth to strike a balance between smoothness and agreement with local characteristics. After applying the LOWESS smoother to the logged series, cases were modeled using ARIMA(1,1,3).

04. ARIMA on the log transformed, LOWESS (quadratic fit) smoothed data. (LOESS)- The `'loess()'` function in R was utilized to implement this approach. Similar to LOWESS, a two-week window was used for smoothing. However, unlike LOWESS, a quadratic polynomial was fitted instead of linear. To mitigate the impact of outliers on the estimation process, a re-descending M estimator with Tukey's bi-weight function was employed. Following the application of the LOESS smoother on the log-transformed data, ARIMA(2,1,2) models were used to analyze cases.

05. ARIMA on the log transformed, SPLINE smoothed data (SPLINE)- In this study, a cubic smoothing spline was utilized for spline smoothing using the `'smooth.spline()'` function in R. The default smoothing parameter, which captures global characteristics, was adjusted to strike a balance between adhering to the time series and achieving suitable smoothness.

For cases, the best degrees of freedom were determined as 40. Following the application of the spline smoother to the log-transformed series, ARIMA(3,1,2) with a drift was employed for modeling COVID-19 cases.

06. ARIMA on log transformed, Running Means data (RM) - A 7-day moving average trend was calculated for the log-transformed data using the 'filter()' function in R., Following this, cases were modeled with ARIMA(3,1,4).
07. Simple Exponential Smoothing on log transformed data (SES) - Simple exponential smoothing was performed using the 'ses()' function from R. The estimated smoothing values for level ('alpha') were 0.34 for cases. The smoothing parameters and initial values were optimized together. The 'initial' argument was set to 'optimal' for selecting initial state values. Forecasting periods ('h') of 1, 7, and 28 were chosen to match the desired time frames.
08. Holt's Exponential Smoothing (without damping the trend) on log transformed data (HOLT)- The 'forecast' package was used to implement the 'holt()' function. The smoothing parameters (α for level, β for trend) were automatically estimated: α was around 0.34 while β was around 0.0004.
09. Holt's Exponential Smoothing (damping the trend) on log transformed data (HOLT (D))– This second model, HOLT(D), is a derivation of the previous method incorporating a damped trend ('damped' set to 'True'). The estimated damping parameter (ϕ) was approximately 0.80.
10. Artificial Neural Networks on the log transformed data (ANN) - The ANN model is chosen as a benchmark for comparison with other models. The 'nnetar()' function is used to identify NNAR (Neural Network Auto Regression) models, which are feed-

forward neural networks with a single hidden layer and lagged inputs for univariate time series forecasting. The fitted models for non-seasonal data are designated as NNAR(p, k), with NNAR(8, 4) for cases in which 20 networks with random starting values are trained, and their forecasts are averaged.

The models were applied to the testing data, and their forecasting performance was evaluated at various time intervals, while assessing the residuals using Median Absolute Deviation (MAD), Root Mean Squared Error (RMSE), and Robust Root Mean Squared Error (RRMSE). When assessing results, more weight was given to MAD and RRMSE since they are less influenced by outliers in the prediction period. The best-performing models were identified based on their scores across multiple evaluation metrics and presented accordingly.

3. RESULTS AND DISCUSSION

All models are compared at different time intervals. This includes one-day, one-week and one-month forecasts with the objective of determining the most effective method among the models for each time frame. The following table consists of the aforementioned results for the COVID-19 cases in Sri Lanka.

Table 1. Forecast Evaluation of Predictions of COVID-19 Cases

	RAW	LOG	LOWESS	LOESS	SPLINE	RM	SES	HOLT	HOLT(D)	ANN
Next Day Predictions										
RMSE	298.9	265.1	243.1	237.3	237.2	218.7	196.6	200.1	196.9	183.9
MAD	103.7	61.8	25.5	16.4	48.7	45.7	42.5	44.1	42.4	38.8
RRMSE	120.2	74.5	24.0	17.2	60.5	59.3	53.5	63.6	54.0	47.1

One week ahead Predictions										
RMSE	306.3	329.3	252.5	277.4	238.4	289.3	257.1	266.5	257.0	272.2
MAD	105.0	117.1	41.5	71.4	45.7	83.3	60.9	73.5	60.9	60.3
RRMSE	125.5	139.3	42.4	77.4	58.1	97.9	76.4	117.9	76.4	72.8
One month ahead Predictions										
RMSE	475.7	516.1	397.5	455.4	280.8	460.8	256.6	328.4	256.5	297.8
MAD	306.8	247.5	246.9	263.4	102.4	235.5	52.9	218.5	52.18	114.5
RRMSE	337.5	296.8	225.6	264.0	111.6	270.5	59.3	278.9	59.1	162.4

Among the models tested, LOESS, LOWESS, smoothing splines, Simple Exponential Smoothing, and the Holt exponential smoothing model with the damped trend were the most reliable for COVID-19 forecasting in Sri Lanka. These models consistently outperformed others across different prediction time frames, especially with regard to the outlier-resistant criteria MAD and RRMSE. For short-term predictions, LOWESS and LOESS models demonstrated exceptional performance. Meanwhile, for mid- to long-term forecasts, the smoothing splines, Simple Exponential Smoothing and Holt exponential smoothing with the damped trend models showed strong predictive capabilities.

4. CONCLUSIONS

This study contributes to the ongoing efforts in predicting the spread of COVID-19 by identifying effective statistical models for forecasting cases in Sri Lanka. The evaluation of various models provides insights into their suitability for different scenarios and future pandemics. The analysis of the brief time series data revealed the presence of trends and random peaks, necessitating the use of techniques such as filtering and smoothing to capture these patterns. Among the models examined, the raw ARIMA model demonstrated

the least effectiveness, struggling to handle the complexities and nonlinear patterns present in real-world data. Traditional smoothing techniques, along with ARIMA, outperformed other filtering and deep learning approaches. Notably, exponential smoothing techniques, smoothing splines, LOWESS, and LOESS models exhibited strong performance, sometimes surpassing even complex machine learning models like ANN. This study emphasizes the efficacy of smoothing and filtering approaches over complex machine learning models. Furthermore, it is important to note that under the availability of large volumes of data, this analysis may be further enhanced by incorporating a validation set into the analysis. However, the careful consideration of the temporal evolution of autocorrelation is necessary to ensure the credibility of the validation set and its appropriateness for evaluating model performance. Finally, further research incorporating external factors and regular updates is recommended to enhance the accuracy of forecasts.

REFERENCES

- Amaratunga, D. Cabrera, J., Diaz-Tena, N., Duan, Y., Ghosh, D., Katehakis, M., Lin, C., Wang, J., and Wang, W. (2023) Adaptive learning methodology with application to forecasting COVID-19 daily cases, submitted for publication.
- Cleveland, W. S. (1979). Robust Locally Weighted Regression and Smoothing Scatterplots. In *Source: Journal of the American Statistical Association*, 74 (368).
- Craven, P., & Wahba, G. (1978). Smoothing noisy data with spline functions. *Numerische Mathematik*, 31(4), 377–403.
- Gecili, E., Ziady, A., & Szczesniak, R. D. (2021, January 7). Forecasting COVID-19 confirmed cases, deaths and recoveries: Revisiting established time series modeling through

novel applications for the USA and Italy.
PLOS ONE, 16(1), e0244173.

Hainz, U., Jenewein, B., Asch, E., Pfeiffer, K. P., Berger, P., & Grubeck-Loebenstien, B. (2005). Insufficient protection for healthy elderly adults by tetanus and TBE vaccines. *Vaccine*, 23(25), 3232–3235.

Martinez, E. Z., Aragon, D. C., & Nunes, A. A. (2020). Short-term forecasting of daily COVID-19 cases in Brazil by using the holt's model. *Revista Da Sociedade Brasileira de Medicina Tropical*, 53, 1–5.

Tawiah, K., Iddrisu, W. A., & Asampana Asosega, K. (2021). Zero-inflated time series modelling of covid-19 deaths in Ghana. *Journal of Environmental and Public Health*, 1–9.