



Modelling COVID-19 incidence in the African sub-region using smooth transition autoregressive model

Eric N. Aidoo¹ · Richard T. Ampofo² · Gaston E. Awashie³ · Simon K. Appiah² · Atinuke O. Adebanji¹

Received: 12 November 2020 / Accepted: 11 February 2021 / Published online: 26 February 2021
© The Author(s), under exclusive licence to Springer Nature Switzerland AG part of Springer Nature 2021

Abstract

Prediction of COVID-19 incidence and transmissibility rates are essential to inform disease control policy and allocation of limited resources (especially to hotspots), and also to prepare towards healthcare facilities demand. This study demonstrates the capabilities of nonlinear smooth transition autoregressive (STAR) model for improved forecasting of COVID-19 incidence in the Africa sub-region were investigated. Data used in the study were daily confirmed new cases of COVID-19 from February 25 to August 31, 2020. The results from the study showed the nonlinear STAR-type model with logistic transition function aptly captured the nonlinear dynamics in the data and provided a better fit for the data than the linear model. The nonlinear STAR-type model further outperformed the linear autoregressive model for predicting both in-sample and out-of-sample incidence.

Keywords COVID-19 · Africa · Nonlinearity · STAR model · Regime switching · Smooth transition

Introduction

The coronavirus disease 2019 (COVID-19) is arguably the world's greatest tragedy of the twenty-first century, which has seen every country across the globe seeking transmission averting strategies and some developed countries producing vaccines in record times (Al-Raei 2021; Bhadra et al. 2020; Le et al. 2020). The COVID-19 is a novel coronavirus caused by severe acute respiratory syndrome coronavirus 2 (SARS-Cov-2) emerged from Wuhan, a capital city of Hubei province in China (Linton et al. 2020; Liu et al. 2020). COVID-19 affects the respiratory system of humans and its infections occur when respiratory droplets of infected

persons are transmitted to susceptible persons in a given population mostly through coughing and sneezing (Linton et al. 2020). This infectious disease has been reported in over 200 countries with more than 25.3 million confirmed cases and 848,000 deaths across the globe as of September 1, 2020 (WHO 2020b). The COVID-19 disease was declared a pandemic by the World Health Organization (WHO) on March 11, 2020 (WHO 2020a).

Within the WHO sub-regions, Africa is ranked fifth with 1,056,120 confirmed cases of coronavirus and 21,999 deaths as of September 1, 2020 (WHO 2020b). The continent of Africa has been predicted to be the next epicentre for COVID-19. In spite of the severe and deadly nature of the disease, there are no vaccine or antiretroviral drugs for treating COVID-19 infections. Different control measures such as physical distancing, ban on social gatherings, quarantine, isolation, and partial/total lockdown of highly infected cities continue to be worldwide. With limited availability of vaccines and the stiff competitions in the acquisition of what is available, the developing world more than ever need enhanced surveillance and intensive-care unit management for infected patients in the fight against the spread of COVID-19. The rapid rate of COVID-19 infections and the emergence of more virulent strains of the virus poses a major global threat but more specifically to the African

✉ Eric N. Aidoo
en.aidoo@yahoo.com

¹ KNUST-Laboratory for Interdisciplinary Statistical Analysis, Department of Statistics and Actuarial Science, Kwame Nkrumah University of Science and Technology, Kumasi, Ghana

² Department of Statistics and Actuarial Science, Kwame Nkrumah University of Science and Technology, Kumasi, Ghana

³ KNUST-Laboratory for Interdisciplinary Statistical Analysis, Department of Mathematics, Kwame Nkrumah University of Science and Technology, Kumasi, Ghana

continent because of the inherent deficiencies in the health-care systems in most of its countries.

The aftermath of increasing number of confirmed cases and associated deaths across the globe, both mathematical and statistical models have been used to describe the transmission dynamics and also forecast the incidence (Maleki et al. 2020; Sanyi et al. 2020; Zhang et al. 2020). The prediction of COVID-19 incidence is of great importance to public health workers and policy-makers to support decision-making regarding distribution of limited resources in areas with high infection rate (Roy et al. 2020a). Public health workers also need prior information on future incidence to prepare for demand of health care facilities.

Existing literature on infectious disease forecasting relies heavily on time-series models (Cortes et al. 2018; Zhang et al. 2014). Within the time-series framework, future incidences of infectious diseases are forecasted based on the historic surveillance data. Linear time-series models such as Autoregressive Integrated Moving Average (ARIMA) models remain the most prominent used statistical technique in COVID-19 incidence forecasting (Ceylan 2020; Maleki et al. 2020; Roy et al. 2020b). However, the use of linear time-series models may not necessarily be appropriate in forecasting future infections, because the dynamics of historical data of COVID-19 is likely to be influenced by social interventions such as ban on social gatherings and partial/total lockdown of cities (Zhang et al. 2020). Such interventions distort the dynamics of data generation process of the disease, which may result in non-constant model parameters. The dynamics of such data may be better described by a nonlinear regime changing time-series models. The regime switching models are nonlinear time-series models that allow model parameters to change for different regimes (Dijk et al. 2002). Among the regime switching models are the Threshold Autoregressive (TAR), Self-Exciting Threshold Autoregressive (SETAR), Markov Switching (MS), and the Smooth Transitional Autoregressive (STAR) models (Dijk et al. 2002; Zivot and Wang 2007).

In this study, the STAR models have been used to model and forecast COVID-19 incidence in African sub-region. The STAR models switch between regimes of the data based on continuous smooth transition functions. The STAR models nest most of the nonlinear times series models such as TAR and SETAR and its forecast performance has been better compared to the nonlinear models (Dijk et al. 2002).

Materials and methods

Data description

The data used in this study were the daily new confirmed cases of COVID-19 across African sub-region. It was

compiled from the official website of World Health Organization (WHO) and span over the period from February 25 to August 31, 2020 (WHO 2020b). An average of 5560 confirmed cases are observed in a day across African sub-region. Figure 1 describes the behaviour of the daily new cases of COVID-19 observed over time in the sub-region. The number of confirmed cases of COVID-19 in the African sub-region showed an upward trend up till around July 23 and declined afterwards.

Smooth transition autoregressive models

The smooth transition autoregressive (STAR) models are nonlinear time-series models, which allow the coefficient of the model to change between different regimes according to a transition function. In this case, we let y_t represents the confirmed new cases of COVID-19 in a particular day t . The two regimes STAR model can be defined for the data as (Teräsvirta 1994):

$$y_t = \mathbf{X}_t \phi^{(1)} (1 - G(y_{t-d}; \gamma, c)) + \mathbf{X}_t \phi^{(2)} G(y_{t-d}; \gamma, c) + \varepsilon_t, \quad (1)$$

where $\mathbf{X}_t = (1, y_{t-1}, \dots, y_{t-p})$, and $\phi_i^{(1)}$ and $\phi_i^{(2)}$ ($i = 0, 1, \dots, p$) represent the model coefficients in both regimes, with p representing the order of the autoregressive process in each regime, $G(y_{t-d}; \gamma, c)$ represents the transition function, and ε_t is the residuals of the model which is assumed to follow a white noise process with zero mean and constant variance σ^2 . The transition variable y_{t-d} is taken to be lagged endogenous variable of y_t with d representing delay parameter ($1 \leq d \leq p$) and c represents a threshold parameter. The transition function $G(y_{t-d}; \gamma, c)$ is a smooth continuous function that controls the switching dynamics, and it is bounded between 0 and 1.

The two most popular choices for the transition function $G(y_{t-d}; \gamma, c)$ are the logistic and exponential functions (Teräsvirta 1994). These two transition functions allow for changes

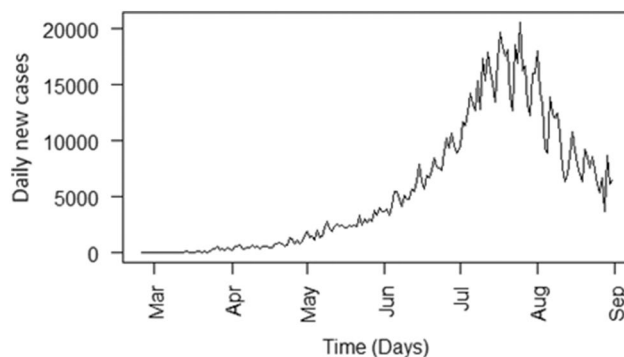


Fig. 1 Temporal pattern of daily new cases of COVID-19 over the study period

in dynamics (Zivot and Wang 2007). The transition function based on the logistic function is defined as:

$$G(y_{t-d}; \gamma, c) = \frac{1}{1 + \exp[-\gamma(y_{t-d} - c)]}, \quad \gamma > 0, \quad (2)$$

where γ determines the speed and smoothness of transition from one regime to another, and the resultant model leads to a logistic STAR (LSTAR) model. When $\gamma \rightarrow \infty$, the LSTAR model becomes a TAR model, and when $\gamma \rightarrow 0$, it becomes a linear AR model. The transition function based on the exponential function is defined as:

$$G(y_{t-d}; \gamma, c) = 1 - \exp[-\gamma(y_{t-d} - c)^2], \quad \gamma > 0, \quad (3)$$

and the resultant model leads to exponential STAR (ESTAR) model. When $\gamma \rightarrow \infty$ or 0, the ESTAR model becomes a linear AR(p) model. The characteristics of the logistic and exponential functions are described in Fig. 2. The logistic function allows different dynamics to occur between the contraction and expansion regimes, whilst exponential function allows similar dynamics to occur with different dynamics in the middle between contraction and expansion period (Terasvirta and Anderson 1992). As γ increase, both logistic and exponential functions become steeper leading to a fast transition between the two regimes.

The model requires that the COVID-19 incidence data be stationary and nonlinear (Teräsvirta 1994). The Zivot–Andrews unit root test (Zivot and Andrews 2002) was adopted for test of stationarity. The null hypothesis of unit root with structural change was tested against an alternative of trend stationary. In addition, the Lagrange multiplier (LM) type test (Teräsvirta 1994) was used for testing nonlinearity in the incidence data. The LM test, the null hypothesis of linearity is tested against the alternative of nonlinear STAR-type model. Assuming d is known, the LM test is equivalent to the test of

$$H_0 : \beta_{1,i} = \beta_{2,i} = \beta_{3,i} = 0, \quad i = 1, 2, \dots, p$$

against the alternative H_1 : ‘ H_0 is not valid’ after performing the auxiliary regression defined as (Luukkonen et al. 1988):

$$y_t = \mu + \sum_{i=1}^p \alpha_{1,i} y_{t-i} + \sum_{i=1}^p \beta_{1,i} y_{t-i} y_{t-d} + \sum_{i=1}^p \beta_{2,i} y_{t-i} y_{t-d}^2 + \sum_{i=1}^p \beta_{3,i} y_{t-i} y_{t-d}^3 + e_t, \quad (4)$$

where e_t is the error term. The LM test statistic asymptotically follows a standard chi-square distribution under the null hypothesis with $3p$ degrees of freedom (Dijk et al. 2002). The value of d has to be specified to perform the LM test. Thus, the test is performed for different values of $d \in (1, 2, 3, \dots, p)$, where d is chosen for which the null hypothesis is rejected. If the null hypothesis is rejected for more than one d , then the appropriate value of d corresponds to the one with smallest P value. The data are considered nonlinear if the null hypothesis is rejected at 5% significance level. After rejecting the null hypothesis of linearity, the appropriate transition function should be selected by testing the following sequence of nested null hypothesis:

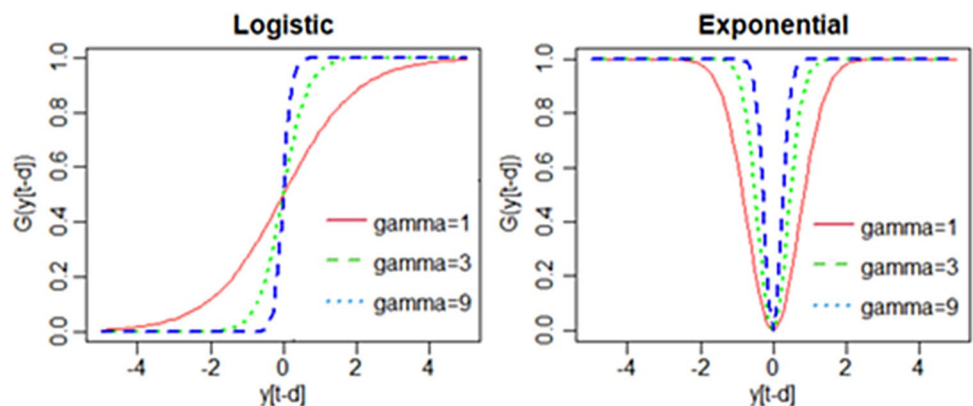
$$H_{0,1} : \beta_{3,i} = 0, \quad i = 1, 2, \dots, p$$

$$H_{0,2} : \beta_{2,i} = 0 \mid \beta_{3,i} = 0, \quad i = 1, 2, \dots, p$$

$$H_{0,3} : \beta_{1,i} = 0 \mid \beta_{2,i} = \beta_{3,i} = 0, \quad i = 1, 2, \dots, p.$$

The decision rule in such a test is that LSTAR should be selected if $H_{0,1}$ is the rejected. Again, ESTAR should be selected if $H_{0,1}$ is not rejected, but $H_{0,2}$ is rejected. On the other hand, LSTAR should be selected if $H_{0,1}$ and $H_{0,2}$ are not rejected, but $H_{0,3}$ is rejected. After identifying a suitable transition function, the coefficients of the selected model can be estimated using nonlinear least-squares method (Teräsvirta 1994; Terasvirta and Anderson 1992). The estimated model may be evaluated by examining the model residuals for the presence of autocorrelation and conditional heteroscedasticity using Ljung–Box (LB) test and autoregressive condition heteroscedasticity (ARCH) test, respectively.

Fig. 2 Characteristics of logistic and exponential function for different values of gamma



The forecast performance of the fitted models may be compared to the linear AR model using mean absolute error (MAE) and root-mean-square error (RMSE). The MAE and RMSE are defined as:

$$MAE = \frac{1}{n} \sum_{t=1}^n |y_t - \hat{y}_t| \tag{5}$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{t=1}^n (y_t - \hat{y}_t)^2}, \tag{6}$$

where y_t represents the observed and \hat{y}_t represents the predicted value, and n is the number of observations.

Results

The analysis was conducted on the log-transformed data to reduce the impact due to high variations in the data. The model building involved 178 observations, starting from February 25 to August 20, 2020, whilst the observations from August 21–31, 2020 were used to validate the forecast performance of the fitted model. The Zivot–Andrews test statistic of -4.644 and the associated critical value of -5.08 at 5% significance level suggest that the data are not stationary. The next step is to test the null hypothesis of linearity as against the alternative of nonlinear STAR model using the LM test. In the LM test, an AR(5) model was used and the order of p was selected to minimize the Akaike information criterion (AIC). The P values of the LM test as shown in Panel A of Table 1 suggest that the null hypothesis of linearity should be rejected at 5% significance level. The results of the test were consistent for all possible values of d . The appropriate value of d was found to be 1, since the p value associated with $d=1$ is minimal compared to the other p values. Having rejected the linearity hypothesis, the appropriate transition function for the STAR model to capture the nonlinear dynamics of the data has to be selected. The results from the sequence of nested hypothesis show that the null hypothesis $H_{0,1}$ is not rejected, but $H_{0,2}$ and $H_{0,3}$

are rejected at 5% significance level (Panel B of Table 1). However, a comparison of the P values associated with $H_{0,2}$ and $H_{0,3}$ suggests that the STAR model with logistic transition function (LSTAR) could be an appropriate model to capture the nonlinear dynamics of the COVID-19 incidence cases. The estimated parameters of the LSTAR model using nonlinear least-squares and the associated goodness-of-fit statistics are presented in Table 2. The estimated LSTAR model was evaluated by examining the models’ residuals for the presence of autocorrelation and conditional heteroscedasticity using LB test and ARCH test, respectively. The p values of the LB test show that there is no serial autocorrelation in the residuals of the model up to lag order p . The

Table 2 Estimated parameters for the LSTAR models and goodness-of-fit statistics

Parameters	Estimate	Standard error	t value	P value
Low regime				
$\phi_0^{(1)}$	0.020	0.053	0.384	0.701
$\phi_1^{(1)}$	-0.584	0.068	-8.591	<0.001
$\phi_2^{(1)}$	-0.765	0.135	-5.686	<0.001
$\phi_3^{(1)}$	-0.315	0.080	-3.932	<0.001
$\phi_4^{(1)}$	-0.335	0.078	-4.291	<0.001
$\phi_5^{(1)}$	-0.381	0.071	-5.365	<0.001
High regime				
$\phi_0^{(2)}$	1.568	0.712	2.201	0.028
$\phi_1^{(2)}$	0.257	0.276	0.931	0.352
$\phi_2^{(2)}$	-0.357	0.372	-0.959	0.338
$\phi_3^{(2)}$	0.602	0.249	2.414	0.016
$\phi_4^{(2)}$	0.635	0.211	3.014	0.003
$\phi_5^{(2)}$	1.697	0.365	4.649	<0.001
γ	6.331	2.034	3.113	0.002
c	0.610	0.096	6.342	<0.001
Model diagnostics				
AIC	-407			
LB(p)	1.682			0.891
ARCH(p)	10.688			0.058

Table 1 P values of the LM test for linearity (Panel A) for different values of d^* and the sequential LM test for appropriate transition function (Panel B)

Null hypothesis	$d=1$	$d=2$	$d=3$	$d=4$	$d=5$
Panel A					
$H_0 : \beta_{1,i} = \beta_{2,i} = \beta_{3,i} = 0$	< 0.001 (85)	< 0.001 (74)	< 0.001 (75)	< 0.001 (49)	< 0.001 (62)
Panel B					
$H_{0,1} : \beta_{3,i} = 0$	0.881				
$H_{0,2} : \beta_{2,i} = 0 \mid \beta_{3,i} = 0$	0.034				
$H_{0,3} : \beta_{1,i} = 0 \mid \beta_{2,i} = \beta_{3,i} = 0$	0.003				

*In parentheses are the test statistic values of the LM test

Table 3 The MAE and RMSE of in-sample and out-of-sample forecasting

Model	In-sample		Out-of-sample	
	MAE	RMSE	MAE	RMSE
Nonlinear LSTAR model	0.208	0.297	0.261	0.366
Linear AR model	0.251	0.385	0.277	0.372

results from the ARCH test show that there is no ARCH effect in the residuals.

The predictive performance of the estimated LSTAR model was further compared with the linear AR(5) model using MAE and RMSE calculated for both models. Smaller values of MAE and RMSE indicate better model performance. The in-sample forecast performance was determined based on the data for the study period (from February 25 to August 20, 2020). The in-sample forecast performance of the nonlinear LSTAR model outperformed the linear AR (5) model (Table 3). The estimated model was used to predict the incidence of COVID-19 beyond the sample period (August 21–31, 2020). The results based on the MAE and RMSE show that the nonlinear model surpasses that of the linear AR (5) model.

Discussion

In this study, the nonlinear smooth transition autoregressive model has been used to model and predict COVID-19 incidence in the African sub-region. The model was based on daily confirmed new cases of COVID-19. Unlike the linear time-series models that assumed constant parameters for the entire data, the nonlinear STAR model allows the model parameters to change for different regimes of the data.

The Zivot–Andrews test showed that the historical incidence data of COVID-19 has a unit root with structural change. The use of Zivot–Andrews test as opposed to the traditional unit root tests was influenced by the fact that, in the presence of nonlinearity, the results from the traditional unit root test are fallacious (Alimi et al. 2017). From the empirical results, the null hypothesis of linearity against the alternative STAR-type nonlinearity was rejected at 5% significance level. In addition, the results of a sequence of nested hypothesis showed that the nonlinear dynamics of COVID-19 incidence is better described by the STAR-type model with logistic transition function (LSTAR model). The presence of such nonlinear behaviour in the daily new cases of COVID-19 may be influenced by the major control strategies such as lockdowns, physical distancing, closure of borders with neighbouring countries, and “stay at home” measures implemented by different countries in the African sub-region (Zhang et al. 2020).

The sign of the estimated parameters of the fitted LSTAR model varies for different regimes. Such behaviour suggests that a model with fixed or constant model parameters will not be appropriate to capture the dynamic patterns of the incidence data. The estimated value of γ indicates that the logistic transition between the two regimes is slow. The forecasting performance of the nonlinear LSTAR model was found to outperform the linear counter AR counterpart based on the MAE and RMSE. The performance of the LSTAR model was consistent in both in-sample and out-of-sample forecast.

conclusion

In conclusion, the case study presented here shows here highlights the need to consider nonlinear dynamics in modelling and predicting the incidence of COVID-19. The use of nonlinear model such as LSTAR can capture any nonlinear dynamics that might be present in the data. The findings of the study could be incorporated into the decision-making process regarding prediction of future incidence of COVID-19.

Author contributions Conceptualization: ENA, RTA, GEA, SKA, and AOA; methodology: ENA and GEA; formal analysis and investigation: ENA, RTA, SKA, and GEA; writing—original draft preparation: ENA, RTA, and GEA; writing—review and editing: SKA and AOA; resources: EN, SKA, and AOA; supervision: SKA and AOA.

Funding No funding was received for conducting this study.

Data availability The dataset used and/or analysed during the current study is available in the public domain.

Code availability Available on request from the corresponding author.

Compliance with ethical standards

Conflict of interest The authors have no conflicts of interest to declare that are relevant to the content of this article.

References

- Al-Raei M (2021) Numerical simulation of the force of infection and the typical times of SARS-CoV-2 disease for different location countries. *Model Earth Syst Environ*. <https://doi.org/10.1007/s40808-020-01075-3>
- Alimi M, Rhif A, Rebai A (2017) Nonlinear dynamic of the renewable energy cycle transition in Tunisia: evidence from smooth transition autoregressive models. *Int J Hydrog Energy* 42:8670–8679. <https://doi.org/10.1016/j.ijhydene.2016.07.131>
- Bhadra A, Mukherjee A, Sarkar K (2020) Impact of population density on Covid-19 infected and mortality rate in India. *Model Earth Syst Environ*. <https://doi.org/10.1007/s40808-020-00984-7>

- Ceylan Z (2020) Estimation of COVID-19 prevalence in Italy, Spain, and France. *Sci Total Environ*. <https://doi.org/10.1016/j.scitotenv.2020.138817>
- Cortes F et al (2018) Time series analysis of dengue surveillance data in two Brazilian cities. *Acta Trop* 182:190–197. <https://doi.org/10.1016/j.actatropica.2018.03.006>
- Dv D, Teräsvirta T, Franses PH (2002) Smooth transition autoregressive models - a survey of recent developments. *Econom Rev* 21:1–47. <https://doi.org/10.1081/etc-120008723>
- Le TT, Andreadakis Z, Kumar A, Román RG, Tollefsen S, Saville M, Mayhew S (2020) The COVID-19 vaccine development landscape. *Nat Rev Drug Discov* 19:305–306
- Linton NM et al (2020) Incubation period and other epidemiological characteristics of 2019 novel coronavirus infections with right truncation: a statistical analysis of publicly available case data. *J Clin Med* 9:538. <https://doi.org/10.3390/jcm9020538>
- Liu J et al (2020) Impact of meteorological factors on the COVID-19 transmission: a multi-city study in China. *Sci Total Environ*. <https://doi.org/10.1016/j.scitotenv.2020.138513>
- Luukkonen R, Saikkonen P, Teräsvirta T (1988) Testing linearity against smooth transition autoregressive models. *Biometrika* 75:491–499. <https://doi.org/10.2307/2336599>
- Maleki M, Mahmoudi MR, Wraith D, Pho K-H (2020) Time series modelling to forecast the confirmed and recovered cases of COVID-19. *Travel Med Infect Dis*. <https://doi.org/10.1016/j.tmaid.2020.101742>
- Roy S, Bhunia GS, Shit PK (2020a) Spatial prediction of COVID-19 epidemic using ARIMA techniques in India. *Model Earth Syst Environ*. <https://doi.org/10.1007/s40808-020-00890-y>
- Roy S, Bhunia GS, Shit PK (2020b) Spatial prediction of COVID-19 epidemic using ARIMA techniques in India. *Model Earth Syst Environ*. <https://doi.org/10.1007/s40808-020-00890-y>
- Sanyi T et al (2020) Analysis of COVID-19 epidemic traced data and stochastic discrete transmission dynamic model. *Sci Sin Math* 50:1070. <https://doi.org/10.1360/ssm-2020-0053>
- Teräsvirta T (1994) Specification, estimation, and evaluation of smooth transition autoregressive models. *J Am Stat Assoc* 89:208–218. <https://doi.org/10.2307/2291217>
- Teräsvirta T, Anderson HM (1992) Characterizing nonlinearities in business cycles using smooth transition autoregressive models. *J Appl Econ* 7:S119–S136. <https://doi.org/10.1002/jae.3950070509>
- WHO (2020a) Coronavirus disease (COVID-19) outbreak situation. World Health Organization. <https://www.who.int/emergencies/diseases/novel-coronavirus-2019>. (Accessed 20 Aug 2020)
- WHO (2020b) WHO coronavirus disease (COVID-19) Dashboard. https://covid19.who.int/?gclid=Cj0KCQjw4f35BRDBARIsAPePBHwxTYv19VVO0wluWvtvktmf2SArOSQHKraGYj-wRXTnIG7_ayeFBMaAgLfEALw_wcB. (Accessed 1 Sept 2020)
- Zhang X, Ma R, Wang L (2020) Predicting turning point, duration and attack rate of COVID-19 outbreaks in major Western countries. *Chaos Soliton Fract*. <https://doi.org/10.1016/j.chaos.2020.109829>
- Zhang X, Zhang T, Young AA, Li X (2014) Applications and comparisons of four time series models in epidemiological surveillance data. *PLoS ONE* 9:e88075. <https://doi.org/10.1371/journal.pone.0088075>
- Zivot E, Andrews DWK (2002) Further evidence on the great crash, the oil-price shock, and the unit-root hypothesis. *J Bus Econ Stat* 20:25–44. <https://doi.org/10.2307/1391541>
- Zivot E, Wang J (2007) *Modeling financial time series with S-plus*. Springer Science and Business Media, Berlin

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