
ORIGINAL ARTICLE

A Vector Autoregression Prediction Model for COVID-19 Outbreak

Qinan Wang¹ | Yaomu Zhou¹ | Xiaofei Chen^{2,3}

¹Lyle School of Engineering, Southern Methodist University, Dallas, TX, 75275, USA

²Department of Statistical Science, Southern Methodist University, Dallas, TX, 75275, USA

³Department of Population & Data Sciences, UT Southwestern Medical Center, Dallas, TX, 75390, USA

Correspondence

Xiaofei Chen, PhD

Department of Statistical Science, Southern Methodist University, Dallas, TX, 75275, USA

Email: xiaofei@smu.edu

Since two people came down a county of north Seattle with positive COVID-19 (coronavirus-19) in 2019, the current total cases in the United States (U.S.) are over 12 million. Predicting the pandemic trend under effective variables is crucial to help find a way to control the epidemic. Based on available literature, we propose a validated Vector Autoregression (VAR) time series model to predict the positive COVID-19 cases. A real data prediction for U.S. is provided based on the U.S. coronavirus data. The key message from our study is that the situation of the pandemic will get worse if there is no effective control.

KEYWORDS

COVID-19, Prediction, Time series data, Vector autoregression, Internal validation

1 | INTRODUCTION

COVID-19 (coronavirus-19) a new type virus, belonging to the Coronaviridae family, spreads from Wuhan, China in 2019.[1] The Coronaviridae family consists of two main subfamilies: Coronavirinae and Torovirinae. These viruses affect the neurological, gastrointestinal, hepatic, and respiratory systems and can be grown by humans, livestock, etc.[2, 3, 4]

Since the appearance, COVID-19 has infected over 59 million people worldwide. [5]The worst situation experienced by the United States (U.S.) followed by the United Kingdom, Italy, France, and Spain. The U.S. has a cumulative 12 million positive cases up to now. It found itself grappling with the worst outbreak after Italy and Spain.[6]The Centers for Disease Control and Prevention (CDC) has verified evidence that COVID-19 is distributed from human to human,

and has also reported that COVID-19 spreads through touching surfaces, close contact, air, or objects that contain viral particles. In the incubation period, it can spread to others. It should be noted that the incubation period and median age of confirmed cases are 3 days and 47 years respectively. [7, 8]

The economic and social disruption caused by the pandemic is devastating. The disease prevention and control is eager for a disease prediction guidance. Efficient models for short-term forecasting has a pivotal role to develop strategic planning methods in the public health system. Under the guidance of the prediction model, we know the severity and the trends of epidemic under different strategies. It can arouse public awareness and help government take the most benefit measures to avoid deaths and reduce infection, such as ordered school closure, case-base measures, the banning of public events, the encouragement of social distancing, and lockdown.

Per literature, a system of differential equations for Susceptible-Infected-Removed (SIR) sequences is a typical mathematical epidemiological model for COVID-19 forecasting.[5, 9, 10, 11, 12, 13, 14] Joining SIR models, Khan et al. proposed the SQUIDER compartmental model to predict the coronavirus 2019 spread [15], and Xu et al. applied the generalized fractional-order SEIR model.[16] The SIR model has a good fitting for the simulation and data of the outbreak in the early stage of the disease. However, the obvious limitations are not limited to that the overall model system has a small external control power, and the number of patients presents a typical exponential growth, which is due to the absence of external drugs and preventive measures.

Other works on COVID-19 prediction has been carried out in Deep Learning and ARIMA (Auto Regressive Integrated Moving Average) univariate time series model. To assess the dynamics of epidemic diseases, time series analysis tools and deep learning are also widely used in publications. Zeroual et al., Shahid et al., and Chimmula et al. performed the Recurrent Neural Network to predict the spread.[17, 18, 19] With time series tools, Alzahrani et al., Sahai et al., and Kumar et al. predicted the COVID-19 by ARIMA univariate model.[20, 21, 22] Deep learning requires a high number of training samples. However, the data we have are still few, so the model generalization is unappealing, namely overfitting. In the time series field, the ARIMA model is quite simple, requiring only endogenous variables and no other exogenous variables. A Stationary is required for the time series, or it is stationary after differencing. Essentially, it lays a shortfall in explaining the causality between different variables.

This article aims to build a generalized VAR (Vector Autoregressive) model for predicting the dynamics of COVID-19 daily cases of the epidemic. VAR is a comprehensive model integrating the advantages of multiple linear regression and the advantages of time series model (the influence of lag term can be analyzed). It applies linear relations to describe a stable system. Under the stationary condition, we can achieve a consistent estimator with the least-square estimation. Besides, VAR can describe the dynamic linear correlation between variables that affect each other, whether used for prediction, interpretation, or sensitivity analysis are clear. With the selected correlated variables among undetected infected, detected deaths, detected recovered, average temperature, precipitation, wind speed, humidity, population density, social trust and civic engagement, that are commonly cited in other epidemiology publications, VAR multivariate model can have a better performance on forecasting and provide an interpretive result. [15, 23, 24, 25, 26, 27, 28]

The correlated variables we choose are useful for analyzing the critical factors driving epidemics. Not only for COVID-19 but may also this model enlighten other epidemics prediction. For the result, some publications tend to be more concerned with the cumulative positive cases, while this article has a very definite awareness of the daily increase cases. A cumulative positive cases prediction is less meaningful than a daily cases increase, since the latter is a better

representative signal for epidemic severity. It is also a critical indicator to access the efficiency of COVID-19 control.

In the next section (Section 2), we describe the data we used in the analyses. The method section (Section 3) elaborates the proposed VAR model and analysis plan. The Results section (Section 4) provides the prediction results by VAR modeling and an internal validation/evaluation of the model. Section 5 discusses the model performance, further improvement, and comparison with other models.

2 | THE DAILY REPORTED COVID-19 DATA

The COVID-19 disease has been reported by CDC (Centers for Disease Control and Prevention) and published in nation and fifty states in the United States by the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University. We obtain the data from <https://covidtracking.com/data/download> maintained by "The Atlantic Monthly Group". The data contain the number of death confirmed, death increasing, death probable, hospitalized, hospitalized cumulative, positive confirmed, positive case viral, positive increasing, etc. The available data is beginning on January 22, 2020 to November 24, 2020 (now).

3 | METHOD

3.1 | A Vector Autoregressive Panel Time Series Model

VAR model was proposed by Christopher Sims in 1980s, using all the current variables in the model to carry out regression for some lagged variables. It is an extension of the AR (autoregression) model, which has been widely used for time series. VAR model takes each endogenous variable as a function of the lag value of all endogenous variables in the system, thus extending the univariate autoregressive model to the "vector" autoregressive model composed of multiple time series variables.

Let \mathbf{X}_t be a causal, stationary multivariate process, then the VAR model can be expressed as:

$$\mathbf{X}_t = \alpha + \Phi_1 \mathbf{X}_{t-1} - \dots - \Phi_p \mathbf{X}_{t-p} + \mathbf{a}_t \quad (1)$$

where $\mathbf{X}_t = (X_{t1}, \dots, X_{tm})^T$ is an $m \times t$ matrix; Φ_k is a real-valued $m \times m$ matrix for each $k = 1, \dots, p$; \mathbf{a}_t is multivariate white noise with covariance matrix $\mathbb{E}[\mathbf{a}_t \mathbf{a}_t^T] = \Gamma_a$; 4. $\alpha = (\mathbf{I} - \Phi_1 - \dots - \Phi_p) \boldsymbol{\mu}$, and $\boldsymbol{\mu} = \mathbb{E}(\mathbf{X}_t)$; $\mathbf{I} = \{1, 1, \dots, 1\}$. Now \mathbf{X}_t is called a VAR (p) process, that is, a vector AR process of order p .

Equation (1) can be expressed in multivariate operator notation way: $\Phi(B) (\mathbf{X}_t - \boldsymbol{\mu}) = \mathbf{a}_t$, where $\Phi(B) = \mathbf{I} - \Phi_1 B - \dots - \Phi_p B^p$ and $B^k \mathbf{X}_t = \mathbf{X}_{t-k}$.

A multivariate process \mathbf{X}_t satisfying the difference equation in Equation (1) is a stationary and causal VAR (p) pro-

cess if and only if the roots of the determinantal equation, $|\Phi(z)| = |I - \Phi_1 z - \dots - \Phi_p z^p| = 0$ lie outside the unit circle. A detailed proof see Brockwell et al. and Reinsel et al. [29, 30]

3.2 | Variables potentially correlated to outcome

Several potential variables might influence the number of COVID-19 positive cases according to literature [15, 23, 24, 25, 26, 27, 28]: undetected infected, detected deaths, detected recovered, average temperature, precipitation, wind speed, humidity, population density, social trust, civic engagement, that are considered in other publications of epidemic prediction.

In Chowdhury et al. [31], climate changes directly affect five infectious disease transmission. Altered climatic conditions may increase the vector biting rate and the vector's reproduction rate and shorten the pathogen incubation period. Furthermore, depending on the report, If the temperature is higher than 25.0°C, there is a significant negative correlation between increasing temperature and pneumonia ($p = 0.017$). [31] That is, if the temperature is decreasing under 25.0°C, pneumonia would spread out faster. [31] In Liu et al. [32], when the temperature is lower than 13.0°C, the number of hospital admission increases, which means the speed of infection also rises up. Those are the reason why COVID-19 positive confirmed cases appear rebound tendency after October. [32]

Depending on data reported by the Tasci et al. [33], during the periods of high, normal, and low humidity, the number of days admitted with pneumonia was higher at high humidity rates ($p < 0.05$). AS a result, the speed of COVID-19 transmission would increase at high humidity situation. In other words, the positive confirmed cases show a significant positive relationship with humidity.[33]

According to Brundage et al. [34], the pneumonia rate has a stronger positive correlation with mortality. The mortality increase would affect COVID-19 spreading out faster than before. However, the number of death increased would happen after COVID-19 transmission rising. [34]

Recovered cases should also have a negative correlation with COVID-19. If recovered cases become more, the number of patients with the virus should be less than before. As a result, fewer patients with the virus would match the lower spread of the virus. When the recovered cases are increasing, the transmission of COVID-19 transmission would be controlled.

3.3 | Model selection

As shown in Section 3.1, we need determine the lag order p of the VAR model. There are diverse criteria, Akaike Information Criterion (AIC), Hannan-Quinn Criterion (HQC), Schwarz Criterion (SC), and Final Prediction Error (FPE), to find the optimal p . Specifically, AIC is an estimator of out-of-sample prediction error and thereby relative quality of statistical models for a given set of data. [35, 36] Suppose that we have a statistical model of some data. Let k be the number of estimated parameters in the model and L be the maximum value of the likelihood function for the model.

Then the AIC value of the model is $AIC = 2k - 2 \ln(L)$. [37, 38] HQC is an alternative to AIC and Bayesian information criterion (BIC). It is given as $HQC = -2L + 2k \ln(\ln(n))$, where n is the number of observations.

Schwarz criterion (SC) is given as $SC = \log(n)k - 2 \log(L(\hat{\theta}))$, where θ is set of all parameter values and $L(\hat{\theta})$ is likelihood of the model returning the data we have, when tested at the maximum likelihood values of θ . Final Prediction Error (FPE) criterion provides a measure of model quality by simulating the situation where the model is tested on a different data set. It is given as $\det \left(\frac{1}{n} \sum_1^n e(t, \hat{\theta}_i) \left(e(t, \hat{\theta}_i) \right)^T \right) \left(\frac{1+d_n}{1-d_n} \right)$, where n is the number of values in the estimation data set, $e(t)$ is a n -by-1 vector of prediction errors, $\hat{\theta}_i$ represents the i -th estimated parameters, d is the number of estimated parameters.

The ordinary least square (OLS) approach is applied to achieve the model estimation. Besides, the model residuals are diagnosed to see if the VAR model assumptions meet.

4 | RESULT

4.1 | Preliminary analysis

According to literature and correlation analysis, we include cumulative death, cumulative recovered patients, temperature and humidity in the VAR model. Considering the positive cases prediction nationwide, we choose the climate data of Washington D.C. that could be representative.

The correlation analyses are shown in Figure 1. Even though 'Death' and 'Humidity' have relatively small correlation coefficients, 0.071 and 0.016, with daily positive case increase, we still keep these two variables. Because the COVID-19 have been verified correlated with cumulative death cases and different humidity.[33, 34]

A descriptive analysis for cumulative death case, cumulative recovered case, temperature, humidity is shown in Figure 2. The cumulative death cases and the cumulative recovered cases presents straight up tendency.

Since co-integration between daily positive cases and other selected variables is required by VAR model. We run the co-integration test (Engle Granger test) on all the variables (series). This test is for daily increase positive cases with other variables. The null hypothesis is that there is no co-integration relationship between the two variables. If the variables are all co-integrated with daily positive cases, we can claim that they have are stably correlated in a long run. Results see Appendix Table 3. As all p -values are less than 0.05, we have all variables co-integrated with daily increase positive cases. Based on the above results, it is considered that there is a stable relationship and there is no spurious regression for the constructed model.

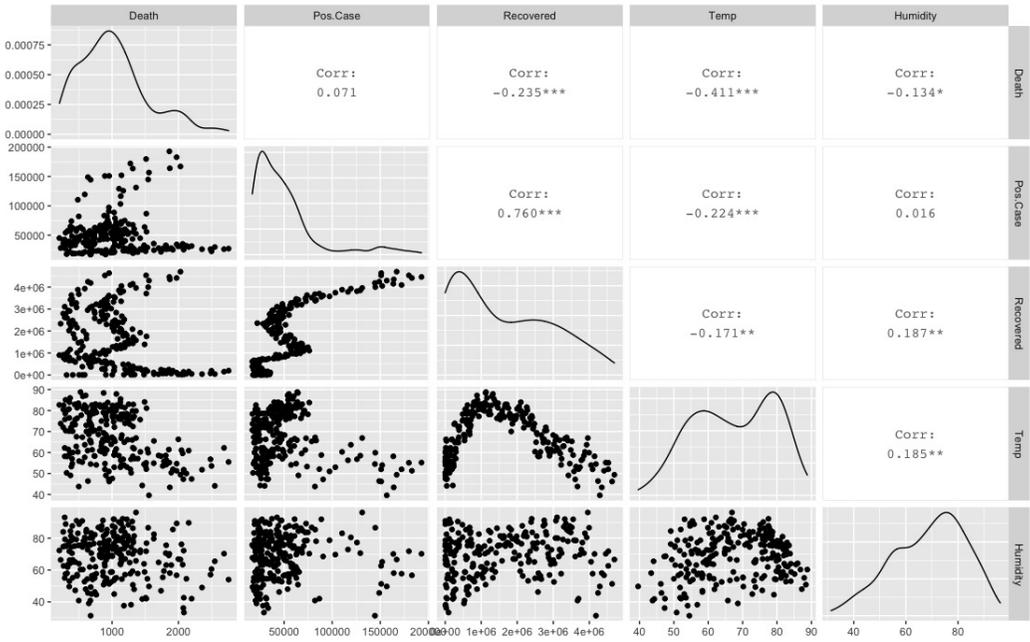


FIGURE 1 Correlation matrix plot. (Death: Cumulative death case; Pos.Case: daily positive case; Recovered: cumulative recovered case; Temp: temperature)

4.2 | Model selection

To determine the lag order of our VAR model, we take AIC, HQC, SC and FPE into consideration (results see Appendix Table 4). The optimal lag order is determined to be 8.

With the suggested lag order 8, we estimate our model using ordinary least square technique. We show the parameter estimates in Table 1.

To verify the assumptions of VAR model, we plot residuals and residuals autocorrelation as shown in Appendix Figure 4. The mean value of residual is almost zero ($-1.45e-14$) and autocorrelation coefficients are within 95% confidence interval (CI; blue dotted line). We also test the residuals by Ljung-Box test and have p-value 0.20 (null hypothesis is that the data are independently distributed). Hence, the fitted model satisfies the assumptions mentioned above: $\mathbb{E}(e_t) = 0$ and $\mathbb{E}(e_t e_{t-k}^T) = 0$, where e_t is the residual at time t .

4.3 | COVID-19 daily positive cases prediction

The objective is predicting the trend of the daily increase positive cases. We predict 30-day daily positive cases starting from July 2, August 21, and November 24, respectively, for internal validation purpose. We pick these three dates for particular reasons. First, 30-day daily increase positive cases after July 2 and August 21 are not fluctuating too

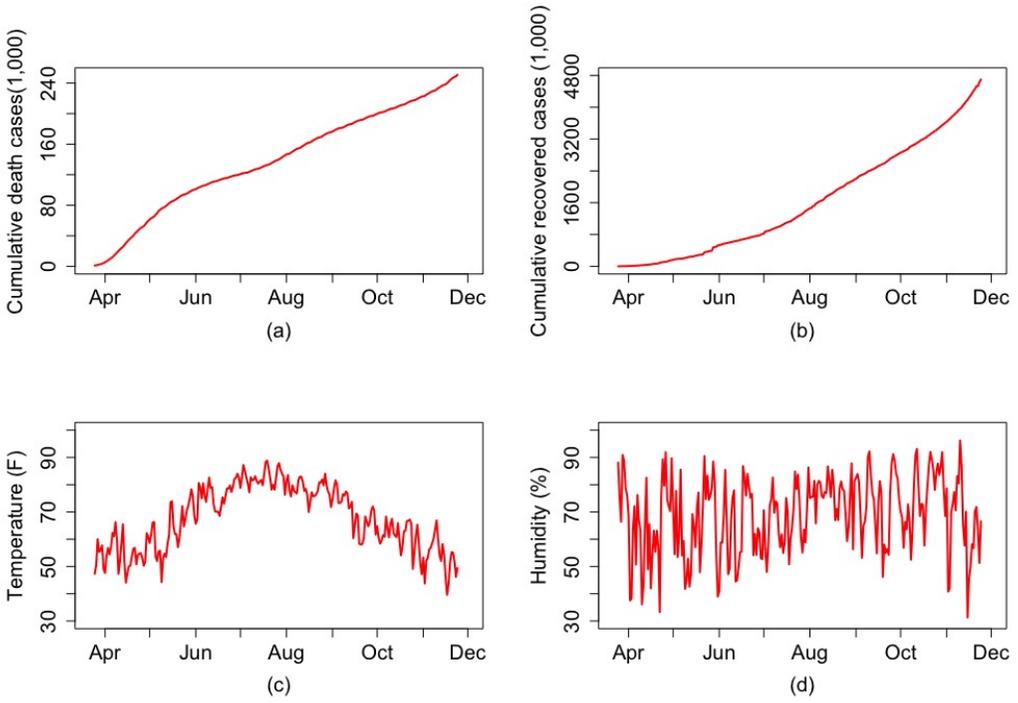


FIGURE 2 Time series from March 25 to November 24, 2020: (a) Cumulative death of COVID-19 cases in US; (b) Cumulative recovered cases; (c) Temperature in Washington D.C.; (d) Humidity in Washington D.C..

TABLE 1 VAR model parameters estimation¹

Lag order/ Variables	Death	Pos.Case	Recovered	Temp	Humidity
1	3.62	2.29	7.49	1.30	4.97
2	-5.07	-1.28	-3.46	-3.68	-1.21
3	3.25	-4.38	-1.53	-2.96	4.12
4	-5.00	-40.62	9.61	2.66	-3.91
5	9.06	2.09	1.67	-1.61	-1.69
6	7.43	-2.46	-2.92	-3.96	8.07
7	-6.89	5.80	-1.50	3.61	-6.47
8	1.53	-1.29	1.94	-5.22	-2.83

Note: ¹ "constant" item is 8327.24.

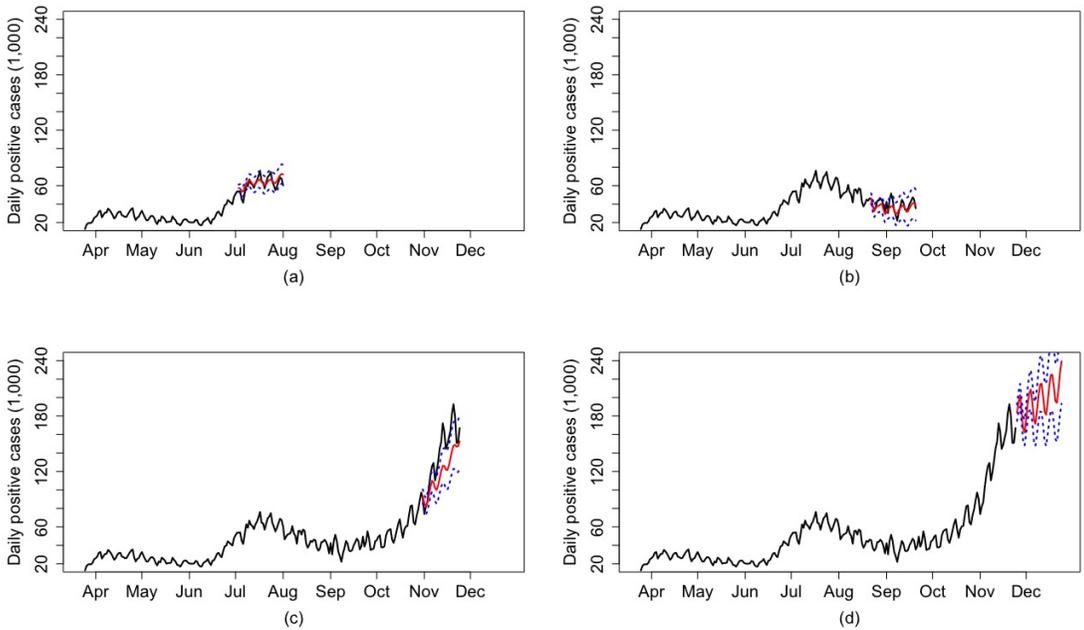


FIGURE 3 Real data (black line) and prediction (red line) with 95% confidence interval (blue dotted line).

(i). For validation purpose: (a) 30 days prediction from July 2, 2020. (b) 30 days prediction from August 21, 2020. (c) 30 days prediction from September 6, 2020.

(ii). Real prediction: (d) 30 days prediction from November 24, 2020.

much. It is a preliminary test on model performance. Critically, the trend after September 6 becomes steep suddenly. However, the trend before September 6 is similar to before July 2 and before August 21. It may be trapped for the model to identify these three conditions. We want to test if the model will predict the correct rapid increase after September 6.

As mentioned, the first three plots (a) – (c) in Figure 3 are for internal validation purpose. As we can see that the model is useful, since the real data (black) is covered by the predicted 95% confidence interval (blue dotted line). To be specific, in Figure 3 (a), the black line is the real selected positive confirmed cases daily data, which presents a stable tendency in the first three months, around 40,000 cases every day. After that, in the middle of June, the real data begin to increase. The short red color line is our prediction using proposed model. The black line and the red line are almost overlapped. One thing need to mention is that the black line fluctuates slightly larger than the red, but the predicted is mostly covered by the 95% confidence interval. It concludes a satisfied prediction. In Figure 3 (b), after the middle of July, the number of daily positive confirmed cases decreases and experiences the first peak of 80,000 cases. However, both the black and red lines show a stable trend then, and the figure shows almost the same appearance as the first 30-day prediction. In Figure 3 (c), the real data experiences a decreasing trend. But, at the beginning of September, the number of daily positive confirmed cases appears a rebound, directed straight up to the second peak value. The peak value even reaches two hundred thousand cases for one day. Our predicted values are a little lower

TABLE 2 The real values and predictions of the daily increase positive cases on Tuesday with 95% confidence interval¹.

Date	Real value	Prediction	Lower 95% CI	Upper 95% CI
2020-07-09	58961	64116	58240	69993
2020-07-16	70446	66497	57820	75173
2020-07-23	71225	66129	56216	76043
2020-07-30	68806	71059	60634	81484
2020-08-25	36588	35466	28817	42115
2020-09-01	42426	30940	21343	40536
2020-09-08	22137	30037	17357	42717
2020-09-15	34904	32092	16668	47515
2020-11-03	86662	86890	77433	96348
2020-11-10	131182	104960	89508	120411
2020-11-17	156722	121455	100033	142878
2020-11-24	167012	152604	123769	181438
2020-12-01		176995	161998	191993
2020-12-08		185548	160309	210787
2020-12-15		194945	159893	229998
2020-12-22		208196	164852	251539

¹ Values on December 1, December 8, December 15, December 22 are blank since real data are not available until now.

than real values. The reason can include Halloween holiday parties and some assemblies because those happened at the end of October, and many COVID-19 cases can be confirmed in early November. Furthermore, those are some extrinsic factor besides the ones in our VAR model. As a result, it is reasonable that the black line is higher than the red prediction line and the confidence interval's upper bound. The success is that the model correctly predicts the rapidly increasing trend after August 21.

The last plot (d) in Figure 3 is our main result that predicts the daily positive COVID-19 cases 30 days later starting from November 24 (now), that is, a prediction for unknown future trend (to December 24). It is obvious that the the future 30-day growth trend will increase if government are not taking any new measures to control the transmission of COVID-19. During the Christmas, the predicted daily positive case is around 240,000 in US.

Table 2 shows a comparison of real values and predictions with a 95% confidence interval. Considering that the daily cases increase data on Monday is partly derived from the cases accumulation over the weekend, we compare the predicted data on each Tuesday with the real values. In Table 2, the real values are generally within the 95% confidence level. For the predictions on November 3, 10, 17, 24, the model predicts 86,890, 104,960, 121,455, 152,604 and upper bounds are 96,348, 120,411, 142,878 and 181,438. The real values exceed upper bounds on November

10 and 17 by around 10,000. It still shows the real values are within the confidence interval since the real values do not deviate from the upper bound too far. The model has good performance of catching a rapid increase trend and regular trend.

5 | DISCUSSION

The study proposed and applied the VAR model for predicting the dynamics of daily COVID-19 positive cases. We selected relevant variables according to literature and checked their correlation coefficients and co-integration. We evaluated our model by comparing the predicted values and real values.

We can introduce more relevant variables in the future to improve the performance if outside force appears to influence viral transmission, control or exacerbate. The most possible variables available may be the estimation of social distance and the number of vaccination. Then the model will be still useful after vaccine comes out. It enables the model to predict the decrease of infections at the vaccination initial stage. It is also the reason why we investigate the application of VAR model on pandemic predictions. The VAR model is different from and better than SIR and ARIMA. Because SIR and ARIMA have an unsatisfied performance when outside force gets involved.

The proposed model can be strongly generalized because it is not limited to specific data, since the structure of the model is constructed. Based on the generalization, this model can be used to predict other epidemics with the same characteristics as COVID-19.

references

- [1] Lu R, Zhao X, Li J, Niu P, Yang B, Wu H, et al. Genomic characterisation and epidemiology of 2019 novel coronavirus: implications for virus origins and receptor binding. *The Lancet* 2020;395(10224):565–574.
- [2] Chen Y, Liu Q, Guo D. Emerging coronaviruses: genome structure, replication, and pathogenesis. *Journal of medical virology* 2020;92(4):418–423.
- [3] Tang B, Wang X, Li Q, Bragazzi NL, Tang S, Xiao Y, et al. Estimation of the transmission risk of the 2019-nCoV and its implication for public health interventions. *Journal of clinical medicine* 2020;9(2):462.
- [4] Wang L, Kraemer R, Borngraeber J. An improved highly-linear low-power down-conversion micromixer for 77 GHz automotive radar in SiGe technology. In: 2006 IEEE MTT-S International Microwave Symposium Digest IEEE; 2006. p. 1834–1837.
- [5] Kufel T, et al. ARIMA-based forecasting of the dynamics of confirmed Covid-19 cases for selected European countries. *Equilibrium Quarterly Journal of Economics and Economic Policy* 2020;15(2):181–204.
- [6] Konarasinghe K. Modeling COVID-19 Epidemic of USA, UK and Russia. *Journal of New Frontiers in Healthcare and Biological Sciences* 2020;1(1):1–14.
- [7] Guan Wj, Ni Zy, Hu Y, Liang Wh, Ou Cq, He Jx, et al. Clinical characteristics of 2019 novel coronavirus infection in China. *MedRxiv* 2020;.
- [8] Maleki M, Mahmoudi MR, Wraith D, Pho KH. Time series modelling to forecast the confirmed and recovered cases of COVID-19. *Travel Medicine and Infectious Disease* 2020;p. 101742.

- [9] Malavika B, Marimuthu S, Joy M, Nadaraj A, Asirvatham ES, Jeyaseelan L. Forecasting COVID-19 epidemic in India and high incidence states using SIR and logistic growth models. *Clinical Epidemiology and Global Health* 2020;.
- [10] Dhanwant JN, Ramanathan V. Forecasting COVID 19 growth in India using Susceptible-Infected-Recovered (SIR) model. *arXiv preprint arXiv:200400696* 2020;.
- [11] Ndiaye BM, Tendeng L, Seck D. Analysis of the COVID-19 pandemic by SIR model and machine learning technics for forecasting. *arXiv preprint arXiv:200401574* 2020;.
- [12] Bastos SB, Cajueiro DO. Modeling and forecasting the Covid-19 pandemic in Brazil. *arXiv preprint arXiv:200314288* 2020;.
- [13] Fanelli D, Piazza F. Analysis and forecast of COVID-19 spreading in China, Italy and France. *Chaos, Solitons & Fractals* 2020;134:109761.
- [14] Chen YC, Lu PE, Chang CS, Liu TH. A Time-dependent SIR model for COVID-19 with undetectable infected persons. *IEEE Transactions on Network Science and Engineering* 2020;.
- [15] Khan Z, Van Bussel F, Hussain F. A predictive model for Covid-19 spread—with application to eight US states and how to end the pandemic. *Epidemiology & Infection* 2020;148.
- [16] Xu C, Yu Y, Yang Q, Lu Z. Forecast analysis of the epidemics trend of COVID-19 in the United States by a generalized fractional-order SEIR model. *arXiv preprint arXiv:200412541* 2020;.
- [17] Zeroual A, Harrou F, Dairi A, Sun Y. Deep learning methods for forecasting COVID-19 time-Series data: A Comparative study. *Chaos, Solitons & Fractals* 2020;140:110121.
- [18] Shahid F, Zameer A, Muneeb M. Predictions for COVID-19 with deep learning models of LSTM, GRU and Bi-LSTM. *Chaos, Solitons & Fractals* 2020;140:110212.
- [19] Chimmula VKR, Zhang L. Time series forecasting of COVID-19 transmission in Canada using LSTM networks. *Chaos, Solitons & Fractals* 2020;p. 109864.
- [20] Alzahrani SI, Aljamaan IA, Al-Fakih EA. Forecasting the spread of the COVID-19 pandemic in Saudi Arabia using ARIMA prediction model under current public health interventions. *Journal of infection and public health* 2020;13(7):914–919.
- [21] Sahai AK, Rath N, Sood V, Singh MP. ARIMA modelling & forecasting of COVID-19 in top five affected countries. *Diabetes & Metabolic Syndrome: Clinical Research & Reviews* 2020;14(5):1419–1427.
- [22] Kumar P, Kalita H, Patairiya S, Sharma YD, Nanda C, Rani M, et al. Forecasting the dynamics of COVID-19 Pandemic in Top 15 countries in April 2020: ARIMA Model with Machine Learning Approach. *medRxiv* 2020;.
- [23] Siedner MJ, Harling G, Reynolds Z, Gilbert RF, Haneuse S, Venkataramani AS, et al. Social distancing to slow the US COVID-19 epidemic: Longitudinal pretest–posttest comparison group study. *PLoS medicine* 2020;17(8):e1003244.
- [24] Behnood A, Golafshani EM, Hosseini SM. Determinants of the infection rate of the COVID-19 in the US using ANFIS and virus optimization algorithm (VOA). *Chaos, Solitons & Fractals* 2020;139:110051.
- [25] Bialek S, Bowen V, Chow N, Curns A, Gierke R, Hall A, et al. Geographic differences in COVID-19 cases, deaths, and incidence—United States, February 12–April 7, 2020 2020;.
- [26] Elgar FJ, Stefaniak A, Wohl MJ. The trouble with trust: Time-series analysis of social capital, income inequality, and COVID-19 deaths in 84 countries. *Social Science & Medicine* 2020;263:113365.
- [27] Bruine de Bruin W. Age differences in COVID-19 risk perceptions and mental health: Evidence from a national US survey conducted in March 2020. *The Journals of Gerontology: Series B* 2020;.

- [28] James N, Menzies M. Cluster-based dual evolution for multivariate time series: Analyzing COVID-19. *Chaos: An Interdisciplinary Journal of Nonlinear Science* 2020;30(6):061108.
- [29] Brockwell PJ, Davis RA, Fienberg SE. *Time series: theory and methods: theory and methods*. Springer Science & Business Media; 1991.
- [30] Reinsel GC. *Elements of multivariate time series analysis*. Springer Science & Business Media; 2003.
- [31] Chowdhury FR, Ibrahim QSU, Bari MS, Alam MJ, Dunachie SJ, Rodriguez-Morales AJ, et al. The association between temperature, rainfall and humidity with common climate-sensitive infectious diseases in Bangladesh. *PLoS One* 2018;13(6):e0199579.
- [32] Liu Y, Kan H, Xu J, Rogers D, Peng L, Ye X, et al. Temporal relationship between hospital admissions for pneumonia and weather conditions in Shanghai, China: a time-series analysis. *BMJ open* 2014;4(7).
- [33] Tasci SS, Kavalci C, Kayipmaz AE. Relationship of meteorological and air pollution parameters with pneumonia in elderly patients. *Emergency Medicine International* 2018;2018.
- [34] Brundage JF, Shanks GD. Deaths from bacterial pneumonia during 1918–19 influenza pandemic. *Emerging infectious diseases* 2008;14(8):1193.
- [35] McElreath R. *Statistical rethinking: A Bayesian course with examples in R and Stan*. CRC press; 2020.
- [36] Taddy M. *Business data science: Combining machine learning and economics to optimize, automate, and accelerate business decisions*. McGraw Hill Professional; 2019.
- [37] Burnham KP, Anderson DR. *A practical information-theoretic approach. Model selection and multimodel inference*, 2nd ed Springer, New York 2002;2.
- [38] Akaike H. A new look at the statistical model identification. *IEEE transactions on automatic control* 1974;19(6):716–723.

Appendix

TABLE 3 Result for Engle-Granger test (co-integration test)

Variables	Statistics	p-value	Co-integration
Cumulative death cases	0.0656	< .0001	Y
Cumulative recovered cases	0.0864	< .0001	Y
Temperature	0.0472	< .0001	Y
Humidity	0.0373	< .0001	Y

TABLE 4 Lag order selection: AIC, HQ, SC, and FPE

Lag order	AIC	HQC	SC	FPE
1	5.57	5.59	5.62	1.57
2	5.50	5.54	5.59	7.82
3	5.48	5.53	5.61	6.52
4	5.47	5.53	5.63	5.67
5	5.45	5.53	5.65	4.76
6	5.44	5.54	5.68	4.45
7	5.42	5.52	5.69	3.42
8	5.39	5.51	5.69	2.49
9	5.39	5.53	5.73	2.55
10	5.39	5.55	5.78	2.76

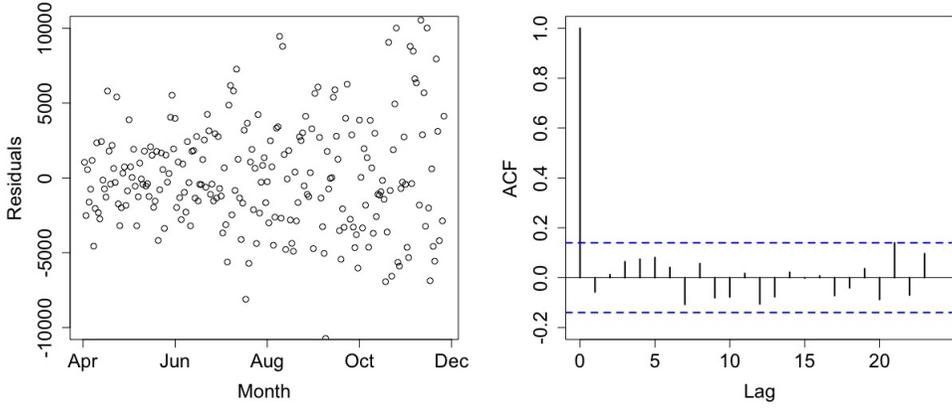


FIGURE 4 Residual plot (left) and its autocorrelation (right). Blue dash lines are upper and lower bound of 95% confidence interval.