

Case Report

Open Access

Time series analysis of Holt model and the ARIMA model facing Covid-19

Xia Jiang¹, Jinming Cao², Bin Zhao^{3*}

¹ Hospital, Hubei University of Technology, Wuhan, Hubei, China.

² School of Information and Mathematics, Yangtze University, Jingzhou, Hubei, China.

³ School of Science, Hubei University of Technology, Wuhan, Hubei, China.

*Corresponding Author: Dr. Bin Zhao, School of Science, Hubei University of Technology, Wuhan, Hubei, China Tel./Fax: +86 130 2851 7572. E-mail: zhaobin835@nwsuaf.edu.cn

Citation: Time series analysis of Holt model and the ARIMA model facing Covid-19: A Case Report. Anna Cas Rep and Ima Surg. 2020; 3(2): 01-09.

Submitted: 12 June 2020; Approved: 15 June 2020; Published: 17 June 2020

ABSTRACT:

Background:

Since the first appearance of the novel coronavirus in Wuhan in December 2019, it has quickly swept the world and become a major security incident facing humanity today. While the novel coronavirus threatens people's lives and safety, the economies of various countries have also been severely damaged. Due to the epidemic, a large number of enterprises have faced closures, employment has become more difficult, and people's lives have been greatly affected. Therefore, to establish a time series model for Hubei Province, where the novel coronavirus first broke out, and the United States, where the epidemic is most severe, to analyze the spreading trend and short-term forecast of the new coronavirus, which will help countries better understand the development trend of the epidemic and make more adequate preparation and timely intervention and treatment to prevent the further spread of the virus.

For the data collected from Hubei Province, including cumulative diagnoses, cumulative deaths, and cumulative cures, we use SPSS to establish the time series model. Since there is no problem of missing data values, we define days as the time variable, remove outliers, and set the width of the confidence interval to 95% for prediction, then use SPSS's expert modeler to find the best-fit model for each sequence. ACF, PACF graphs of the residuals, and Q-tests are used to determine whether the residuals are white noise sequences and to check whether the model is a suitable model. Holt model is used for the cumulative number of diagnoses, and ARIMA (1,2,0) model is used for cumulative cures and deaths. Similarly, we also collect data for the US, including the cumulative number of diagnoses, cumulative deaths, and cumulative cures. For the three groups mentioned above, ARIMA (2,2,6) model, ARIMA (0,2,0) model, and ARIMA (0,2,1) model are used respectively.

Findings:

From our modeling of the data, the time series diagrams of the real the fitted data almost overlap, so the fitting effect of the Holt model and the ARIMA model we use is very suitable. We compare the predicted values with the real values of the same period and find that the epidemic situation in Hubei Province has basically ended after May, but the epidemic situation in the United States has become more severe after May, so the Holt model and the ARIMA model are also very appropriate in predicting the epidemic situation in short-term. Interpretation:

Because the Chinese government has always put the safety of people 's lives in the first place, when the epidemic broke out, it decisively closed the city of Hubei Province.

One side is in trouble, all sides support, they concentrate all resources of whole country to save Hubei Province at the expense of the economy only in order to save more people. Now we can clearly see that the epidemic has been controlled in China and the whole country is developing in a good direction. The situation in the United States, on the other hand, is also influenced by the social environment.

Keywords: COVID-19; time series analysis; Holt model; ARIMA model; SPSS

1. Introduction

Since December 2019, many cases of cough, dyspnea, fatigue and fever have been reported in Wuhan, China, in a short period of time. Unexplained cases of pneumonia with major symptoms. The Chinese health authorities and the Centre for Disease Control and Prevention(CDC) promptly identified the causative agent of these cases as a novel coronavirus, and the World Health Organization (WHO) in November 2009 identified the virus as the main cause of pneumonia. On January 10, 2020, it was named COVID-19. on January 22, 2020, the People's Republic of China the State Council Information Office held a press conference on pneumonia to prevent and control neo-coronavirus infections. On the same day, the Center for Disease Control and Prevention of the People's Republic of China released a plan for the prevention and control of neo-coronavirus infection of pneumonia. Includes COVID-19 epidemiological study, specimen collection and testing, follow-up and management of close contacts and information, education and training. Risk communication with the public.

And now in May in our country is almost under control, but in the meantime the novel coronavirus has brought the life of the people and the development of the national economy. Coronaviruses are a large family of viruses that are known to cause influenza as well as Middle East Respiratory Syndrome (MERS) and Severe Acute Respiratory Syndrome (SARS). Coronaviruses are a large family of viruses known to cause influenza, as well as Middle East Respiratory Syndrome (MERS) and Severe Acute Respiratory Syndrome (SARS). More serious diseases such as respiratory syndrome (SARS). Novel coronaviruses are new strains of coronaviruses that have never been found in humans before. Common signs of coronavirus infection in humans include respiratory symptoms, fever, cough, shortness of breath, and difficulty breathing. In more severe cases,

the infection can lead to pneumonia, severe acute respiratory syndrome, kidney failure, and even death. Not only that, but in later cases we also found a proportion of asymptomatic infections, which is defined as the absence of relevant clinical symptoms (e.g., self-perceived or clinically recognizable symptoms and signs such as fever, cough, sore throat, etc.), respiratory and other specimens of coronavirus etiology (usually referred to as a nucleic acid test) or a positive serum-specific immunoglobulin M (IgM) antibody test. An asymptomatic infected person is not a confirmed case and therefore is not counted as a newly infected patient. announced by the NCHS on March 4, 2020. This provision was not changed in the seventh edition. Asymptomatic infected persons actually comprise two segments of the population: the first is those who are insidiously infected, with no or very mild symptoms. The other part of the population is in a latent phase after infection and may develop symptoms in the future.

Asymptomatic infections can be divided into two scenarios: first, an infected person tests positive for nucleic acid, and after a 14-day incubation period, none of them have any Self-perceptible or clinically recognizable signs and symptoms, always in an asymptomatic state of infection; secondly, positive nucleic acid tests in infected persons. The absence of any self-perceived or clinically recognizable signs and symptoms at the time of sampling, but the subsequent appearance of some clinical manifestation, i.e., latent "Asymptomatic infection" status. Unfortunately, there is currently no specific treatment for diseases caused by novel coronaviruses. However, many of the symptoms can be managed and therefore need to be treated according to the patient's clinical condition. In addition, supportive care for infected patients may be very effective.

By establishing COVID-19 spread models between China and the United States, we can clearly see the different results of

different policies and measures on the control of the epidemic. This is also the significance of our model. Through accurate data and rational

analysis, we can provide the future development trend and control strategy research for the world in the fight against the epidemic².

2. Methods

Data

The data of Hubei Province is derived from the Health Commission of Hubei Province on its official platform from 20 January, 2020. The Hubei Province data collected in this paper includes cumulative deaths, cumulative cures and cumulative diagnoses from 20 January, 2020 to 28 April, 2020.³

The data of the United States comes from the domestic data platform, the Real-time Big Data Report of the Epidemic. The data collected in this paper includes cumulative deaths, cumulative cures and cumulative number of diagnoses from 29 February, 2020 to 28 April, 2020.⁴ The model

Through the collected data, we conduct a time series analysis of the novel coronavirus^{5,6,7}. Because there is no data missing, we import the data into SPSS, define the day as the time variable, remove the outliers, and make time series graphs. The most suitable fitting models are automatically found by the expert modeler, which include the cumulative deaths, cumulative cures and cumulative diagnoses. The explanation of abbreviations used in time series analysis are listed below. **Table 1: abbreviations mentioned in time series analysis**

abbreviations	Explanations for each abbreviation				
	Auto Regressive. The AR model is a statistical a method for processing time series that predicts the current period of x_t using the previous periods of the same variable, e.g., $x_t x_t$ to $x_{t,t}$, performance, and assuming they are a linear relationship. since this is a development from linear regression in regression				
AR(<i>p</i>)	analysis, except that instead of using x to predict y , x predicts x (itself).				
	Moving Average. The MA model is one of the model parametric spectral analysis methods and is a commonly used model in modern spectral estimation. q -order moving average model (MA(q)) autocorrelation				
MA(q)	coefficients have q-order truncation.				
	Autoregressive Integrated Moving Average. The ARIMA model, also known as the differential integrated moving average autoregressive model, is one of the time series prediction analysis methods. In ARIMA (p , d , q), AR is "auto regressive", p is the number of autoregressive terms; MA is "sliding average", q is the number of autoregressive terms. The number of sliding mean terms, d				
	the number of times (in steps) the difference was made to make it asmooth sequence.				
ARIMA(<i>p</i> , <i>d</i> , <i>q</i>)					
ACF	Autocorrelation Function. Autocorrelation, also called serial correlation, is the correlation of a signal with itself at different points in time. Informally, it is a function of the similarity between two observations on the time				
	difference between them.				
PACF	Partial Autocorrelation Function. Rather than finding the correlation between a lag like <i>ACF</i> and the current, Partial Autocorrelation Function finds the correlation between the residual and the next lag value. Thus, if there is any hidden information in the residuals that can be modeled by the next lag, we may get a good cor- relation, and We will use the next lag as a feature when				
	modeling.				
TS	Time Series				
Q test	It is a method to determine whether an outlier in a set of data is suspicious or not.				
UCL	Upper Control Limit. Upper bound of the confidence interval.				
LCL	Lower Control Limit. Lower bound of the confidence interval.				

As the Table 1 shows, the ACF and the PACF are used to check whether the model is fitted. The ACF is the autocorrelation coefficient and does not control for other variables. The partial autocorrelation coefficient PACF, on the other hand, is the autocorrelation coefficient calculated after controlling for other variables, since it hacks out the effects of other variables. ACF is mainly used for MA model, and PACF is mainly used for AR model. The Table 2 listed below illustrates the meaning of ACF and PACF plots in AR model, MA model, ARMA model.

Table 2: The meaning of ACF and the PACF plots inAR(p) model, MA(q) model and ARMA (p, q) model

models	ACF	PACF
AR(p)	Gradual decay, i.e. trailing	Post p-order truncation
MA(q)	Post q-order trun- cation.	Gradual decay, i.e. trailing
ARMA(p, q)	Gradual decay, i.e. trailing	Gradual decay, i.e. trailing

TS Model-based method for estimation in Hubei Province

Based on the given data of cumulative number of diagnoses in Hubei Province, we use the expert modeler to process the data and obtain the Holt model8,9 to describe it. The corresponding equation set shown below.

> $St = \alpha Xt + (1 t \alpha) (Stt_1 + Ttt_1)$ $Tt = \beta (St t Stt_1) + (1 t \beta) Ttt1$ $X^{t}(N) = St + NTt$

The related mathematical symbols used above are listed in the following Table 3

 Table 3: Mathematical symbols used in equation set

Symbols	Meanings for each symbol
t	Current period
X _t	Actual observations in period t
S _t	Estimated level at period t
T_t	Predicted trend at period t
Х _t N	Estimated value before period <i>m</i>
α	Horizontal smoothing parameter
β	Trend smoothing parameter
m	The number of periods ahead of the forecast, i.e., the forecast step size

Finally, we find that $\alpha = 1$, $\beta = 0.1$.

In order to determine that the selected Holt model can correctly describe the cumulative number of diagnoses, we use white noise10 for residual test. From Figure 1, the ACF and PACF graphs of the residuals can be seen that the autocorrelation coefficients and partial correlation coefficients of all lag orders are not significantly different from 0.From Figure 2, it can also be seen that the P value obtained by the Q test is 1, that is, we cannot reject the null hypothesis and think that the residual is a white noise sequence, so the Holt model can describe the cumulative number of diagnoses well.

Fig.1: The ACF and PACF graphs of the residuals in cumulative number of diagnoses in Hubei Province case



Fig.2: The Q test on the residual of the case of cumulative number of diagnoses in Hubei Province case

		Model Fit statistics			Ljung-Box Q(18)				
Model	Number of Predictors	Stationary R- squared	R-squared	Normalized BIC	Statistics	DF	Sig.	Number of Outliers	
The cumulative confirmed	0	.423	.994	15.254	2.273	16	1.000	0	

From the analysis above, we can conclude that the Holt model can describe the cumulative number of diagnoses well.

Similarly, the expert modeler evaluates that the cumulative number of cures conforms the ARIMA (1, 2, 0) model11. The related equations are demonstrated followed.

 $(^{1 t \Sigma P} i=1^{\alpha i Li}) (1t L)^{dyt} = \alpha 0 + (^{1 + \Sigma q} i=1^{\beta i Li})^{st}$

 $(1 t \alpha_1) (1 t L)^{2yt} = \alpha 0 + st$

The related mathematical symbols used above are listed in the following Table 4.

Т	able	4:	Mathematica	l s	vmbol	s used	in	equation	IS
					,				

Symbols	Meanings for each symbol
р	Number of autoregressive items
q	Moving average number of items
d	The number of differential steps made to become a smooth sequence
L	Lag operator
ε _t	Because $\{\epsilon_t\}$ is a white noise sequence, $E\{\epsilon_t\}=0$
Ρ 1 τΣ α, L ⁱ i=1	AR(<i>p</i>) model
$ \begin{array}{c} q \\ 1 + \Sigma p_i L^i \\ i=1 \end{array} $	MA(q) model
1 t L ²	2nd order difference
y _t	Number of people on day t
at	Coefficients of the i-th order lagged term in the $AR(p)$ model
b _t	Coefficients of the i-th order lagged term in the MA(<i>a</i>) model

The lag operator used in equation ,2,3 indicates the difference. The Table 5 below shows the specific functions of lag operator. Table 5: Functions of lag operator

Functions	Corresponding equations				
<i>d</i> -order dif- ferential	$\Delta^d y_t = (1 - L)^d y_t$				
Seasonal differ- ence (<i>m</i> for pe- riod)	$\Delta y_t - \Delta y_{t-m} = (1-L) (1-L^m) y_t$				

We also use white noise to test this model for residuals. As Figure 3 below shows, the ACF and PACF graphs of the residuals can be seen that the autocorrelation coefficients and partial correlation coefficients of all lag orders are not significantly different from 0. Therefore, the ARIMA (1, 2, 0) model can well describe the cumulative number of cured people.

Fig.3: The ACF and PACF graphs of the residuals in cumulative number of cures in Hubei Province case



Similar to the cumulative number of people cured in Hubei Province, the cumulative number of deaths in Hubei Province also conforms to the ARIMA (1, 2, 0) model. The corresponding equation and the related symbol meanings are the same as equation.

Next, we use white noise for the residual test. From Figure 4, the ACF and PACF graphs of the residuals can be seen that the autocorrelation coefficients and partial correlation coefficients of all lag orders are not significantly different from 0. we can find that the ARIMA (1, 2, 0) model can also describe the cumulative death toll well.

Fig.4: The ACF and PACF graphs of the residuals in cumulative number of deaths in Hubei Province case



TS Model-based method for estimation in the United States

Based on the given data of America, we use the expert modeler to process the data and we find that all of the cumulative number of diagnoses, deaths and cures conform the ARIMA model. However, the parameters setting of each group of them are not identical.

After processing, it is found that the cumulative diagnoses in the United States applies the ARIMA (2, 2, 6) model. The corresponding equation is the same with equation .

Next, we performed a residual test on the model based on white noise. As can be seen from Figure 5, the ACF and PACF graphs of the residuals, the autocorrelation coefficients and partial correlation coefficients of all lag orders are not significantly different from 0. From Figure 6, it can also be seen that the P value obtained from the Q test of the residual is 0.304, that is, we cannot reject the null hypothesis, and think that the residual is a white noise sequence.

Therefore, the ARIMA (2,2,6) model can well describe the cumulative number of diagnoses.



Fig.6: The Q test on the residual of the case of cumulative number of diagnoses in U.S. case

		M	odel Fit statisti	cs.	Lju	ng-Box Q(18	1)	
Model	Number of Predictors	Stationary R- squared	R-squared	Normalized BIC	Statistics	DF	Sig.	Number of Outliers
The cumulative confirmed	0	.793	.999	19.470	17.247	15	.304	

The cumulative number of cures conforms ARIMA (0, 2, 0) model, which is equal to 2nd order difference equation. The related equation set is similar to the equation

As usual, we should use white noise to perform a residual test. As the Figure 7 shows below, the ACF and PACF graphs of the residuals can be seen that the autocorrelation coefficients and partial correlation coefficients of all lag orders are not significantly different from 0. Fig.7: The ACF and PACF graphs of the residuals in cumulative number of cures in U.S. case



At last, we use expert modeler to process the data of cumulative deaths in U.S., then it is found that it conforms the ARIMA (0, 2, 1) model. The related equation set still conforms with the equation

Then we perform a white noise residual test. As can be seen from Figure 8, the ACF and PACF graphs of the residuals, the autocorrelation coefficients and partial correlation coefficients of all lag orders are not significantly different from 0.

Fig.8: The ACF and PACF graphs of the residuals in cumulative number of deaths in U.S.



3.Results

The TS Model-based method results in Hubei Province

We set the width of the confidence interval to 95%, then use the Holt model and ARIMA model to fit and predict the cumulative number of people diagnosed, cumulatively cured and cumulatively died in Hubei Province respectively. The obtained results shown in the following figures.







Date







Date

As can be seen in the Figure 9, 10, 11, the time series plots of the real and fitted data almost overlap, and the Holt and ARIMA models fit the original data well.

At the same time, after 28 April, the epidemic situation in Hubei Province has been controlled, and the cumulative number of diagnoses will basically not increase dramatically. The Holt model and ARIMA model can also well predict the cumulative diagnoses, cumulative cures and cumulative deaths. The Table 3 listed below contains the predicted values of cumulative number of diagnoses, cures and deaths in Day 101 to 105 by using the equation a n d the Figure 9, 10, 11.

Table 3: short-term predicted values in Hubei Prov-ince

Day	Cumulative diagnoses	Cumulative cures	Cumulative deaths
101	68140	63616	4512
102	68152	63616	4512
103	68164	63616	4512
104	68176	63616	4512
105	68188	63616	4512

The TS Model-based method results in the United States

Our analysis of the U.S. epidemic situation also set the width of the confidence interval to 95%. The results obtained by fitting and predicting the number of cumulative diagnoses, cumulative cures and cumulative deaths using the ARIMA model are shown in the following figures.





It can be seen from Figures 12, 13, and 14, the time series plots of the real and fitted data almost overlap, and the Holt and ARIMA models fit the original data well.

At the same time, after 28 April, the cumulative number of diagnoses, cumulative deaths and cumulative cures will continue to increase substantially in the short term, which is related to the policies adopted by the US government to combat the epidemic.

This shows that the ARIMA model can also predict the cumulative diagnoses, cumulative deaths and cumulative cures.

The specific predicted values of cumulative number of diagnoses, cures and deaths in Day 61 to 65 are listed in the Table 4 below. Table 4: short-term predicted values in U.S.

Day	Cumulative diagnoses	Cumulative cures	Cumulative deaths
61	1049197	151006	60312
62	1051063	161873	62006
63	1080926	172741	63722
64	1122025	183609	65458
65	1169266	194477	67217

4. Discussion

At the time of the virus outbreak, everyone is living in panic, fearing for their lives and the lives of their families and the safety of the country.

We use SPSS12 to accurately get the models we need, such as Holt model and ARI-MA model, and then use these models to fit the sequences, and estimate the model parameters based on the sequence values. Finally, we perform residual tests on the models with white noise to check whether the model is applicable.

Time series analysis of COVID-19 gives the course, direction and trend of the epidemic and predicts the likely development of future epidemics to what state. This will give us some guidance in our lives, such as in the response to the epidemic what measures should be taken to intervene in the development of the epidemic, to save more lives.

We can see from the results that not only the epidemic situation of China and the United States selected period is different, but also in epidemic development after 60 days or so. It can be seen that outbreak in the United States will continue to present an exponential growth, but the epidemic situation of Hubei province is basically controlled, probably around 100 days in Hubei province in May outbreak was completely under control. The outcome of the epidemic is different, which is largely related to the measures taken by the country to combat the epidemic, as can be seen from the attitude of China and the United States towards the epidemic13. That's why it's important to have the right understanding in the face of an epidemic, and not to be blindly arrogant and underestimate the seriousness and potential dangers

of an epidemic.

Time series can not only be used in the analysis of infectious diseases, but also in many social disciplines such as measurement¹⁴ and economy¹⁵. The data order and size in time series contain the information of the objective world and its changes, and represent the dynamic process. Therefore, the main purpose of time series analysis is to understand the considered dynamic system, predict future events, and control future events through intervention^{16, 17}.

5. Limitations

When establishing the model, we regard the data with large fluctuations as outliers. In fact, there are many more complex models that can catch these outliers. At the same time, when we make predictions, the ARIMA model is only suitable for short-term prediction. Over a certain period of time, the predicted value will not change any more, due to the theory of the model. So when solving this problem, we can assume the predicted values as the observed values, thus the long-term prediction would be possible, but the difference between the truly observed data might be larger and larger.

Since the epidemic is only predicted in the short term, it can be seen from the analysis chart of epidemic situation in the United States that the cumulative number of people who are cured, died and diagnosed, cases are all moving in an increasing direction. However, in practice, the number of people in these categories should reach a stable value in the end.

Conflict of interest

We have no conflict of interests to disclose and the manuscript has been read and approved by all named authors.

Acknowledgement

This work was supported by the Philosophical and Social Sciences Research Project of Hubei Education Department (19Y049), and the Staring Research Foundation for the Ph.D. of Hubei University of Technology (BSQD2019054), Hubei Province, China.

References

1. https://baike.baidu.com/item/2019%E6% 96%B0%E5%9E%8B%E5%86%A0%E7%8A%B6 %E7%97%85%E6%AF%92/24267858?fr=aladdin

2. http://www.xinhuanet.com/2020-02/01/c_1125518723.htm

3. http://wjw.hubei.gov.cn/

4. https://voice.baidu.com/act/ newpneumonia/newpneumonia/?city=%E7%BE%8E%E5%9B%BD-

%E7%BE%8E%E5%9B%BD

5. Zhou Xuan. Epidemiological characteristics and time series analysis of hand, foot and mouth disease in Wuzhou City from 2014 to 2018[D]. Nanning:Guangxi Medical University.2019. (in Chinese)

6. HE Xiaonan, SONG Xiaohui, ZHU Xin. Prediction of hand-foot-mouth disease incidence in Luoyang based on ARIMA model[J]. Modern Preventive Medicine.2019(03). (in Chinese)

7. Yang Xiaodi. Epidemiological characteristics and time series analysis of hemorrhagic fever in nephrotic syndrome in Jilin province[D]. Changchun:-Jilin University.2019. (in Chinese)

8. Shang Wenli, Zhang Liting, Li Shichao, et al. Dynamic prediction of oil well production based on Holt exponential smoothing model[J]. Automation and Instrumentation, 2018, 033(004):68-70. (in Chinese)

9. https://baike.baidu.com/item/ Holt%E6%A8%A1%E5%9E%8B/22192556?fr=aladdin

10. https://baike.baidu. com/item/%E7%99%B-D%E5%99%AA%E9%9F%B3/10280741?fr=aladdin

11. https://baike.baidu.com/item/ ARIMA%E6%A8%A1%E5%9E%8B/10611682?fr=aladdin

12. http://www.fx361.com/ page/2020/0409/6542570.shtml

13. https://baike.baidu.com/item/ spss/2351375?fr=aladdin

14. Sun Tonghe. Application of time series analysis in the field of measurement[J]. Mapping and Spatial Geographic Information, 2013, 036(003):12-13.

15. Gu Lan.The application of time series analysis to the economy [M]. Beijing: China Statistics Press, 1994. (in Chinese)

16. Zhou Yongdao. Time series analysis and applications [M]. Beijing: Higher Education Press, 2015. (in Chinese)

17. He Shuyuan. Applied time series analysis [M]. Beijing: Peking University Press, 2007. (in Chinese)