

Forecasting Covid-19 Confirmed Cases Deaths And recovered using Vector Autoregression

Dr.N.Poompavai

Department of computer science.

*Bishop packiam Arokiaswamy college of Education,
Pudukkottai -622303.*

Email Poompavai84@gmail.com

Abstract

This article is about the recent complications faced by the people and government of the entire world. The complication is all about an infectious disease caused by a newly discovered virus called coronavirus (Covid-19). This has caused heavy damage to people in the society and also to the government. It took over lakhs of people's lives and dreams. In such cases, the government needs to take precaution. They should be able to predict the future using the statistical report. Here in this article, we have used a Vector Autoregression algorithm for time series data to predict the future, which shows how many people of India are going to be affected, recovered and etc in the next five days. We have done this using the statistical report of India. We have proved this with good accuracy. This article is highly useful to predict the future and using this, government can take precaution before the outbreak of a pandemic may cause damage or loss to the society.

Keywords: *Infectious disease, Coronavirus, Precaution, Pandemic*

1.Introduction

First of all, let us know what Vector Auto Regression is. Vector Auto Regression (VAR) which is mentioned as VAR in many cases is a multivariate forecasting algorithm and is a scholastic process model used to capture the linear among the multiple time series variables. We call it as autoregressive because each variable is represented as a consequence of a past value, that is, the predictors are time delayed values of the series. This is used when one wants to predict multiple series variables. This was developed by Sims in the year 1980. This algorithm is used in many industry, supply, economic, constructions, consumption, saving and etc. So, if we use an autoregressive algorithm in cases of Covid-19, it will be very appropriate. Here in this article, we have collected the data of Covid-19 in India from the website [kaggle.com](https://www.kaggle.com). So now, by implementing the algorithm in the collected data, we will reach an accuracy which may help us to predict the future, that is, we are going to predict the number of fatality, recovery, and also the new cases to be affected due to Covid-19.

2.Vector autoregression:

Vector autoregression (VAR) model is an augmentation of univariate autoregression model to multiple time series data. Here all the variables are dependent on itself, that is, they are treated as endogenous. This is the most successful, flexible and easy model used for the analysis of multivariate time series data. It often provides superior prediction to those from univariate time series. In addition to data illustration and forecasting, this model is also used for analysis and structural intervention. The structure of this model allows us to test restriction around multiple equations and it also tests whether the coefficients on all regression of the lag are zero. This corresponds to testing the null value in the lag order.

Before adding the variables to the VAR model we should recheck it carefully because adding unrelated variables to the VAR model reduces accuracy by increasing the estimation errors. The typical AR model equation is given below

Table 1. Data Format

Date	Confirmed	Deaths	Recovered
2020-04-21	20080	645	3975
2020-04-22	21370	681	4370
2020-04-23	23077	721	5012
2020-04-24	24530	780	5498
2020-04-25	26283	825	5939

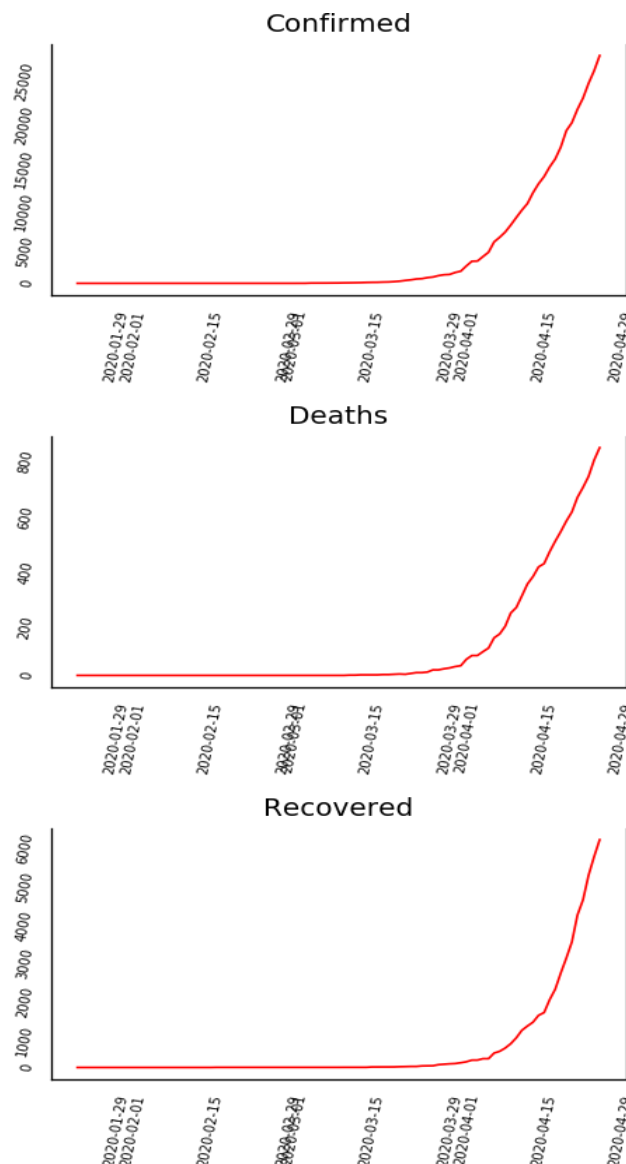


Figure 1. Graphical Representation of Covid Data

3.Experimentalresults:

This section is intended to discuss the results which is generated using the methods given in section 4.1 and 4.2.

3.1 Granger's causality test:

This is a statistical hypothesis test in deciding whether one time series is useful in forecasting the other. This is a technique that seeks the direction of causality between the import and export files. According to the Granger causality test, if a signal in X_1 "Granger-causes" (g-causes) signals X_2 then

the past values of X_1 should contain information that help to predict X_2 . This technique was developed by Granger in 1960. Its formula is based on the linear regression model of scholarly process. The illustration of formula is given below, Consider a bivariate linear autoregressive models whose variables are X_1 and X_2 :

$$X_1(t) = \sum_{j=1}^p A_{11,j} X_1(t-j) + \sum_{j=1}^p A_{12,j} X_2(t-j) + E_1(1)$$

$$X_2(t) = \sum_{j=1}^p A_{21,j} X_1(t-j) + \sum_{j=1}^p A_{22,j} X_2(t-j) + E_2(2)$$

Here,

- a. p is the maximum number of lagged observation included in the model
- b. Matrix A contains the coefficients of the model
- c. X_1 and X_2 are variables
- d. E_1 and E_2 are residuals

Where if the variance of E_1 is reduced by the insertion of the X_2 terms in the first equation then it is said that X_2 g-causes X_1 . This can also be said as X_2 g-causes X_1 if the coefficients in matrix A_{12} are different from zero.

We have tested the Granger's causality for the data collected by us and the resulted output is given in Table 2

Table 2. Granger's Casuality Matrix

	Confirmed_x	Deaths_x	Recovered_x	Confirmed_y	Deaths_y	Recovered_y
Confirmed_x	1.0	0.0	0.0	0.0	1.0	0.0
Deaths_x	0.0	1.0	0.0	0.0	1.0	0.0
Recovered_x	0.0	0.0	1.0	0.0	0.0	1.0
Confirmed_y	0.0	0.0	0.0	1.0	0.0	0.0
Deaths_y	0.0	0.0	0.0	0.0	1.0	0.0
Recovered_y	0.0	0.0	0.0	0.0	0.0	1.0

3.2 Cointegration Test:

Cointegration test is a statistical property of a time series variable. To know more about cointegration first let us know what is ‘order of integration’. The order of integration means the differentiation done to convert the non-stationary time series to stationary. Now, in multiple series there exists a linear combination of multiple series that has order of integration less than that of individual series, and hence the collection of such series is called cointegration. When multivariate series is cointegrated that means they will have a long run significantly.

Now by implementing the formula of cointegration by Johansen in the data collected we will get the result as shown in Table 3.

Table 3. Johansen cointegration test

Name	Teststat	C(95%)	Signif	Confirmed
129.4	24.2761	True		
Deaths	52.74	12.3212	True	
Recovered	9.21	4.1296	True	

3.3 Splitting The Series Into Test And Train Data

To reduce the complication the series is split into train and test data. Here we have splitted the data from 2020-01-22 to 2020-04-20 as train data and data from 2020-04-21 to 2020-04-25 as test data. After splitting the series into train and test data is fitted and the algorithm is implemented. Here we have divided the series into train and test. We have got 90 training data and 5 testing data. And therefore the shape of the train and test data are (90, 3) and (5, 3).

3.4 Test For Stationary To Make The Series Stationary:

The time series we have collected should be stationary in case of the VAR model as it is going to be forecasted. It is complicated to check all the series according to stationarity. So, to reduce this complication we have this technique of test for stationary. Now to check stationarity without complication we have used a popular method namely Augmented Dickey-Fuller Test (ADF), these methods are also called unit root tests or suite of tests.

If a series is found non-stationary we make use of the suite of tests. ADF tests the null value in the time series. This is the large and complicated set of time series models. The ADF test statistic uses the negative number.

Formula of ADF test is given in Equation 4

$$Oy_t = \alpha + \beta t + \gamma y_{t-1} + \delta_1 Oy_{t-1} + \dots + \delta_{p-1} Oy_{t-p+1} + \epsilon t \quad (3)$$

Here,

a. $y(t)$ = variable

b. $y(t-1)$ = lag1 of time service

c. $\Delta y(t-1)$ = first difference of the time series at time (t-1)

d. β = coefficient of the time trend

e. α = constant

f. p = lag order of autoregressive process

$g.Y$ = nullhypothesis

Implementing ADF test in the collected data we will get the result shown in Table 4, Table 5, Table6 respectively.

Table 4. Augmented Dickey-Fuller Test for Confirmed.

Augmented Dickey-Fuller Test on "Confirmed"
Null Hypothesis: Data has unit root. Non-Stationary.
Significance Level = 0.05
Test Statistic = 2.7015
No. Lags Chosen = 11
Critical value 1% = -3.517
Critical value 5% = -2.899
Critical value 10% = -2.587
=> P-Value = 0.9991. Weak evidence to reject the Null Hypothesis.
=> Series is Non-Stationary.

Table 5. Augmented Dickey-Fuller Test for Death

Augmented Dickey-Fuller Test on "Deaths"
Null Hypothesis: Data has unit root. Non-Stationary.
Significance Level = 0.05
Test Statistic = -2.5725
No. Lags Chosen = 12
Critical value 1% = -3.518
Critical value 5% = -2.9
Critical value 10% = -2.587
=> P-Value = 0.0988. Weak evidence to reject the Null Hypothesis.
=> Series is Non-Stationary.

Table 6. Augmented Dickey-Fuller Test for Recovered

Augmented Dickey-Fuller Test on "Recovered"
Null Hypothesis: Data has unit root. Non-Stationary.
Significance Level = 0.05
Test Statistic = 8.2093
No. Lags Chosen = 9
Critical value 1% = -3.515
Critical value 5% = -2.898
Critical value 10% = -2.586
=> P-Value = 1.0. Weak evidence to reject the Null Hypothesis.
=> Series is Non-Stationary.

So this shows us that none of our time series is stationary so to convert these non-stationary series to stationary we need to differentiate the series.

On differentiating the series for the first time we got all the output as non-stationary so now we have two choices either to pursue with 1st difference or we need to differentiate the series one or more times until the series may get converted to stationary. On differentiating the series for the fifth time we will get the following output,

Table 7. Augmented Dickey-Fuller Test on "Confirmed" after fifth difference

Augmented Dickey-Fuller Test on "Confirmed"
Null Hypothesis: Data has unit root. Non-Stationary.
Significance Level = 0.05
Test Statistic = -4.3026
No. Lags Chosen = 12
Critical value 1% = -3.525
Critical value 5% = -2.903
Critical value 10% = -2.589
=> P-Value = 0.0004. Rejecting Null Hypothesis
=> Series is Stationary.

Table 8. Augmented Dickey-Fuller Test on "Death" after fifth difference

Augmented Dickey-Fuller Test on "Deaths"
Null Hypothesis: Data has unit root. Non-Stationary.
Significance Level = 0.05
Test Statistic = -10.1675
No. Lags Chosen = 8
Critical value 1% = -3.519
Critical value 5% = -2.9
Critical value 10% = -2.587
=> P-Value = 0.0. Rejecting Null Hypothesis
=> Series is Stationary.

Table 9. Augmented Dickey-Fuller Test on "Recovered" after fifth difference

Augmented Dickey-Fuller Test on "Recovered"
Null Hypothesis: Data has unit root. Non-Stationary.
Significance Level = 0.05
Test Statistic = -10.6075
No. Lags Chosen = 12
Critical value 1% = -3.521
Critical value 5% = -2.903
Critical value 10% = -2.589
=> P-Value = 0.0. Rejecting Null Hypothesis
=> Series is Stationary.

And the above stated output is the resultant output of the series in differentiating it for the fifth time.

3.5 Selection Of The Order [P] VarModel :

Now, we need to select the right order of the VAR model. We fit increasing orders of VAR models and pick any one order which gives a model with least Akaike information criterion (AIC). We can also check other best fitting comparison estimations and they are Bayesian Information Criterion (BIC), Final Prediction Error (FPE) and Hannan-Quinn Information Criterion (HQIC).

Given below is the resultant output we got by fitting the correct order,

Table 10. Best Fitting Comparison Estimations

	AIC	BIC	FPE	HQIC
0	22.03	22.12	3.7e+09	22.07
1	21.19	21.55	1.5e+09	21.33
2	20.37	21.00	7.0e+08	20.62
3	20.14	21.04	5.1e+08	20.50
4	19.48	20.65	2.9e+08	19.95
5	18.21	19.65	8.2e+07	18.78
6	15.83	17.54	7.6e+06	16.51
7	13.61	15.59	8.4e+05	14.40
8	12.59	14.84	3.1e+05	13.49
9	11.56	14.08	1.1e+05	12.57

4. Check For Serial Correlation Of Residual [Errors] Using Durbin Watson Statistic:

This is used to check if there are any leftover patterns in the errors. That means there are some patterns that are still left to be explained by the model. In such cases, we have three choices that can be done either we can increase the order of the model or we can induce more predictors or we can look for differential algorithm to model the time series. To check correlation of errors we can use Durbin Watson Statistics.

Formula we use to check the correlation is mentioned below,

$$DW = \frac{\sum_{t=1}^T ((e_t - e_{t-1})^2)}{2 \sum_{t=1}^T e_t^2}$$

The resultant value from the implemented formula may vary between 0 to 4. If the value is closer to 0 that means it contains positive correlation, if the value is closer to 4 that means it has negative correlation and if the value is closer to 2 that means there is no serial correlation in the series.

Now let us implement the formula in the data we have collected,

Confirmed : 2.18

Deaths : 2.17

Recovered : 2.34

The above mentioned values are the resultant output we got from implementing the formula Durbin Watson Statistic. So, now the output value is closer to 2 and so we may not have any significant correlation in the multivariate time series data we have collected.

4.1 Forecasting Var Model:

In order to forecast, the model expects the lag order number of observations from the past data. This is because the model lags various time series in the dataset. So, we need to provide many previous values as indicated by the lag order used by the model. The fitted lag order of the collected data is given below,

```
array([ [-1388., -175., -1244.],
       [ 2024., 128., 1232.],
       [-1652., -34., -930.],
       [-1008., -18., 631.] ])
```

The above described data has been forecasted but it is on the scale of training data used by the model. To bring it back to its original form we need to de-differentiate it as many times we have

differentiated the original input data. Here we differentiated the original input data 5 times, so now we need to de-differentiate this resulted forecasted data 5 times.

4.2 Inverting The Transformation To Get The Real Forecast:

Now, let us de-differentiate the forecasted data to get the original data. The de-differentiation of the forecasted data is given below,

Table 11. Forecasting Covid19 dataset

Date	Confirmedforecast	Deathsforecast	Recoveredforecast
2020-04-21	20246	635	3685
2020-04-22	21232	674	4249
2020-04-23	22873	708	4695
2020-04-24	24451	762	5303
2020-04-25	26299	799	5744

And finally the original input data is resulted by de-differentiation which is mentioned above.

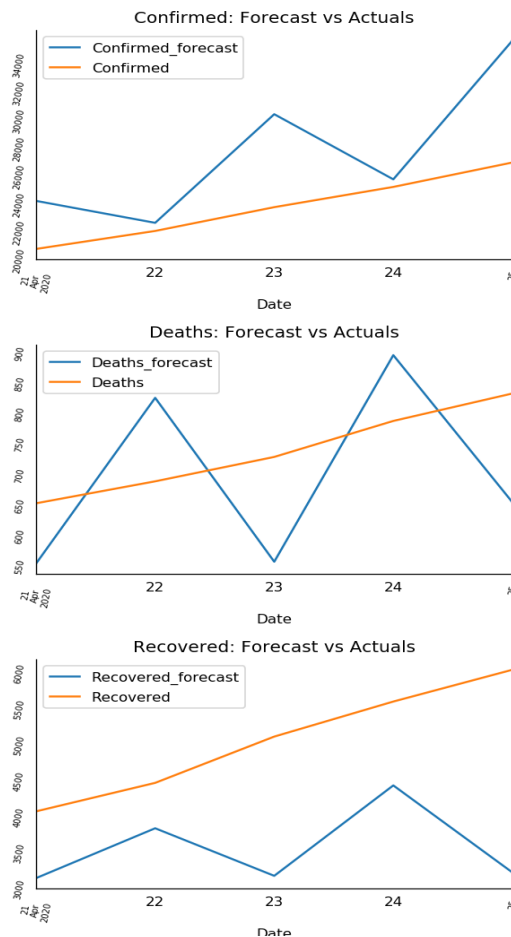


Figure 2. Plotting ForecastVs Actual

5. Evaluating forecast:

To evaluate the forecast we need to compute a comprehensive set of metrics. The metric we have used hereto evaluate the accuracy is MAE (Mean Absolute Error). MAE is the simplest error metric in regression that is understandable. On evaluating the forecast in the collected data according to MAE we will get the following output,

Forecast Accuracy of: Confirmed Cases mae : 120.5973

Forecast Accuracy of: Deaths mae : 14.2689

Forecast Accuracy of: Recovered mae : 223.0366

6. Conclusion:

In this article we have studied about the efficiency vector auto regression from its scratch to the end detail. Using the collected data of covid-19 pandemic in India we have evaluated the successful model using Vector Auto Regression algorithm with good accuracy. So, now we can predict the future cases, deaths and recovery using this model which helps government to take necessary actions in precautionary.

References

- [1] Andrews, D. W. K. (1999), “Consistent Moment Selection Procedures for Generalized Method of Moments Estimation,” *Econometrica*, 67, 543–564. DOI: 10.1111/1468-0262.00036. [Crossref], [Web of Science[®]], [Google Scholar]
- [2] Bekaert, G., Engstrom, E., and Ermolov, A. (2017) “Macro Risks and the Term Structure of Interest Rates,” *Finance and Economics Discussion Series* 2017-058. Washington, DC: Board of Governors of the Federal Reserve System. DOI: 10.17016/FEDS.2017.058. [Crossref], [Google Scholar]
- [3] Bernanke, B. S., and Mihov, I. (1995) “Measuring Monetary Policy,” NBER Working Paper No. 5145. [Google Scholar]
- [4] Brüggemann, R., Jentsch, C., and Trenkler, C. (2016), “Inference in VARs with Conditional Heteroskedasticity of Unknown Form,” *Journal of Econometrics*, 191, 69–85. DOI: 10.1016/j.jeconom.2015.10.004. [Crossref], [Web of Science[®]], [Google Scholar]
- [5] Castelnovo, E. (2016), “Monetary Policy Shocks and Cholesky VARs: An Assessment for the Euro Area,” *Empirical Economics*, 50, 383–414. DOI: 10.1007/s00181-015-0930-2. [Crossref], [Web of Science[®]], [Google Scholar]
- [6] Chaussé, P. (2010), “Computing Generalized Method of Moments and Generalized Empirical Likelihood with R,” *Journal of Statistical Software*, 34, 1–35. DOI: 10.18637/jss.v034.i11. [Crossref], [Web of Science[®]], [Google Scholar]
- [7] Donovan, P., and Hall, A. R. (2018), “The Asymptotic Properties of GMM and Indirect Inference under Second Order Identification,” *Journal of Econometrics*, 205, 76–111. DOI: 10.1016/j.jeconom.2018.03.006 [Crossref], [Web of Science[®]], [Google Scholar]
- [8] Gospodinov, N. (2010), “Inference in Nearly Non-stationary SVAR Models with Long-Run Identifying Restrictions,” *Journal of Business & Economic Statistics*, 28, 1–12. DOI: 10.1198/jbes.2009.08116. [Taylor & Francis Online], [Web of Science[®]], [Google Scholar]
- [8] Gospodinov, N., Kan, R., and Robotti, C. (2014), “Spurious Inference in Unidentified Asset Pricing Models,” *Federal Reserve Bank of Atlanta Working Paper Series*, 2014-12. [Google Scholar]
- [9] Gospodinov, N., Kan, R., and Robotti, C. (2017), “Spurious Inference in Reduced-Rank Asset Pricing Models,” *Econometrica*, 85, 613–1628. DOI: 10.3982/ECTA13750. [Crossref], [Web of Science[®]], [Google Scholar]
- [10] Gospodinov, N., and Ng, S. (2015), “Minimum Distance Estimation of Possibly Noninvertible Moving Average Models,” *Journal of Business & Economic Statistics*, 33, 403–

417. DOI: 10.1080/07350015.201 [Taylor & Francis Online], [Web of Science[®]], [Google Scholar]
- [11] Gouriéroux, C., Monfort, A., and Renne, J.-P. (2017), “Statistical Inference for Independent Component Analysis: Application to Structural VAR Models,” *Journal of Econometrics*, 196, 111–126. DOI: 10.1016/j.jeconom.2016.09.007. [Crossref], [Web of Science[®]], [Google Scholar]
- [12] Guay, A., and Normandin, M. (2018), “Identification of Structural Vector Autoregressions Through Higher Unconditional Moments,” Unpublished Working Paper. [Google Scholar]
- [13] Hall, A. R. (2005), *Generalized Method of Moments*, New York: Oxford University Press. [Google Scholar]
- [14] Hall, A. R., and Inoue, A. (2003), “The Large Sample Behavior of the Generalized Method of Moments Estimator in Misspecified Models,” *Journal of Econometrics*, 114, 361–394. DOI: 10.1016/S030476(03)00089 [Crossref], [Web of Science[®]], [Google Scholar]
- [15] Hall, A. R., Inoue, A., Jana, K., and Sin, C. (2007), “Information in Generalized Method of Moments Estimation and Entropy-Based Moment Selection,” *Journal of Econometrics*, 138, 488–512. DOI: 10.1016/j.jeconom.2006.05.006. [Crossref], [Web of Science[®]], [Google Scholar]
- [16] Hall, A. R., and Peixe, F. P. M. (2003), “A Consistent Method for the Selection of Relevant Instruments,” *Econometric Reviews*, 7, 269–288. DOI: 10.1081/ETC-120024752. [Taylor & Francis Online], [Google Scholar]
- [17] Hansen, L. P. (1982), “Large Sample Properties of Generalized Method of Moments Estimators,” *Econometrica*, 50, 1029–1054. DOI: 10.2307/1912775. [Crossref], [Web of Science[®]], [Google Scholar]
- [18] Hansen, L. P., Heaton, J., and Yaron, A. (1996), “Finite Sample Properties of Some Alternative GMM Estimators Obtained from Financial Market Data,” *Journal of Business and Economic Statistics*, 14, 262–280. DOI: 10.2307/1392442. [Taylor & Francis Online], [Web of Science[®]], [Google Scholar]
- [19] Herwartz, H. (2015), “Structural VAR Modelling with Independent Innovations—An Analysis of Macroeconomic Dynamics in the Euro Area Based on a Novel Identification Scheme,” Working Paper, University of Göttingen. [Google Scholar]