# Time Series Augmentation based on Vector Auto Regression and Long Short Term Memory method for Air Quality Prediction

**Udaya Bharathi Rambha [1]\***          **Maruvada Seshashayee [1]**

*[1] GITAM (Deemed to be University), Visakhapatnam, India*

**Abstract:** Air quality prediction has gained much attention in recent years due to the interest of the people and the government. Air quality prediction helps to take necessary action to improve the air quality and public health. Variousexisting methods involve applying the air quality prediction based on the deep learning and Autoregressive Integrated Moving Average (ARIMA) model. Existing approaches have the limitation of vanishing gradient problems and unstable performance in prediction. In this research, the hybrid method of Vector Auto Regression (VAR) and Long Short Term Memory (LSTM) method is proposed to improve the performance of the Air Quality Index. The Indian Air Quality data are collected from the publicly available Central Pollution Control Board (CPCB). The LSTMmodel is suitable for analyzing the time series prediction due to its capacity to process the data in sequence and ability tostore essential features for the long term. The VAR model performs the normalization based on multivariate data characteristics and augments the data to make the data suitable for the LSTM training. The proposed VAR - LSTM model has an RMSE of 2.633 value; existing SARIMA and ARIMA have 3.932 and 2.896 RMSE, respectively.

**Keywords:** Air Quality Prediction, Deep Learning, Long Short Term Memory, Normalization, and Vector Auto Regression.

## 1.     Introduction

Recently, Air quality has attracted much attentionfrom the people and government due to its profound effect on health and the environment. Monitoring stations are assigned in a city to monitor air quality and other factors. People and the government require air quality to improve air quality by controlling pollution. Air quality prediction is challenging as it depends on several factors, such as air quality spatial-temporal dependencies and weather patterns [1]. Pollutant compositions and concentrations are complicated functions in outdoor air pollutants, and exposure has adverse health impacts. Air pollutants of central outdoor in cities areVolatile Organic Compounds (VOCs), Nitrogen

Oxides (NOx), Carbon Monoxide (CO), Sulfur dioxide (SO2), Particle Matter (PM), Ozone (O3),

pesticides, and metals [2]. Atmospheric PM long-term Exposure decreases lung function and premature death [3]. Data-driven models have been followed inmuch research to find air quality to understand thevarious input parameters of statistical correlation.Machine learning and deep learning were recently applied for forecasting 40 air quality levels [4]. The deep learning method effectively models PM2.5-time series and complex meteorological features [5].

The existing methods apply historical-based prediction in previous research results using simple regression methods, Extreme Learning Machine (ELM), machine learning-based solution, and LSTMor some neural networks for prognosis [6]. Deep learning is popular due to its powerful non-linear fitting ability, and it helps to integrate the different airquality monitor stations' information [7].

Deep learning techniques provide the effective performance to model and predict air quality. Currentair quality prediction research applies deep learning to capture air quality accurately for spatiotemporal patterns and measure the impact of air quality challenges. Existing methods mainly use RNN-based models to predict air quality of temporal dependencies. However, these methods have a vanishing gradient problem that doesn't support a large receptive field, which renders them unable to the long-term prediction of ideal results [8 – 10].

The proposed solution overcomes this problem by augmenting the input data with VAR and improves the learning rate of LSTM, which results in a better prediction of AQI.

1.     The hybrid method VAR and LSTM were applied for Air Quality Prediction to improvethe model's efficiency. The VAR model performs augmentation and normalization for time series data.

2.     Normalized data applied to the LSTM modelimproves the prediction performance for the AirQuality Prediction process.

3. The VAR-LSTM model has a lower error ratethan the existing methods in Air Quality prediction due to the capacity to augment the data.

This paper is organized as a literature review in Section 2, and the proposed method is providedin Section 3. The result is given in Section 4, and the conclusion of this research paper is given in Section 5.

## 2. Literature Review

Air quality prediction is an essential model for air pollution prevention and management. Recently, various models were proposed for air qualityprediction to improve performance efficiency. Some of the notable research in air quality prediction were reviewed in this section.

Ma [11] proposed a Transferred Learning Bidirectional Long Short Term Memory (TL-BiLSTM)model for air quality prediction. The BiLSTM modelhelps learn the long-term dependences and transferthe knowledge from smaller temporal resolution to larger temporal resolution. Various rolling windows were used in the developed method to evaluate the model performance. The developed method of air quality prediction performance was tested in the casestudy of Guangdong, China. The developed method has a lower error rate than existing machine-learning methods in air quality prediction. The vanishing gradient problem in the developed model affects the performance of the prediction process.

Ma [12] proposed a Transfer Learning-based Stacked BiLSTM model to predict the air quality of a new station. The transfer learning and deep learningmethod are combined to perform the prediction thatknowledge gathered from the existing station. Theproposed TLS-BiLSTM model is tested in a casestudy in Anhui, China. The Rolling window wasapplied in the developed method for a training sequence. This method pre-trains the model based onsource data and fine-tunes the model to performprediction. The transfer learning method learns thenumerical patterns in surrounding air quality stationsand learns knowledge to build a prediction model.The developed method has a lower prediction errorthan the existing method in air quality prediction, andthe proposed method was applied to the missing data.

Jin [13] proposed a hybrid deep learning predictor based on Convolutional Neural Network (CNN) forlong-term air quality prediction. Empirical ModeDecomposition (EMD) was   applied  to decompose the input data, and the CNN modelclassified components into a fixed number of groupsbased on frequency characteristics. Each group wastrained with Gated Recurrent Unit (GRU) network asasub-predictor, and three GRU model results were usedfor prediction results. The developed method wastested on Beijing data in air quality prediction. Thedeveloped method has a limitation of overfitting problem that affects the performance of the model.

Krishan [14] applied the LSTM model to predict the air quality index in Delhi, India. Factorssuch as pollutant level, traffic data, meteorological conditions, and vehicular emission were used to predict. The prediction was performed on hourly concentration on 2008-2010 data to analyze the model's performance. The developed model hasmany performances in the air quality prediction, and the LSTM model has the limitation of vanishing gradient problems in the prediction.

Abhilash [15] applied the ARIMA model to predict the air quality index in Bengaluru. The data from 14 pollutant monitoring stations from 2013to 2016 were used to test the model's performance. ARIMA model has considerable performancein air quality prediction. The feature selection performance in the ARIMA model was low, and the model's error rate was high.

Yonkang [16] introduced a hybrid deep learning model that embraces the merits of the stationarywavelet transform (SWT) and the Nested Long ShortTerm Memory (NLSTM)to improve the prediction quality in the problem of hour-ahead air quality forecasting. A framework that leverages several NLSTM recurrent neural networks is constructed to output forecasting results for different sub-signals, respectively. The main limitation of this work is thatit only performs the prediction for PM2.5 values in the given dataset. Although the forecasting process and results for other air-quality indices are similar tothe PM2.5, the deep learning technique actually can perform transit learning from one index to another.

Janarthanan [17] employed the combination of Support Vector Regression (SVR) and Long Short-Term Memory (LSTM) based deep learning model isused to classify the AQI values. The deep learning mechanism accurately predicts the AQI values and helps plan the metropolitan city for sustainable development. The expected AQI value can control pollution by incorporating road traffic signal coordination, encouraging people to use public transportation, and planting more trees in some locations. This work has included many climatic parameters to predict AQI, and it is suitable for only aparticular metropolitan city.
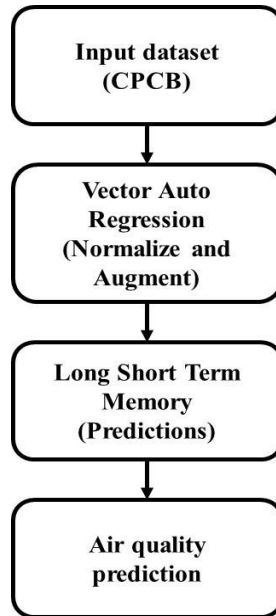
## 3.    Proposed Method



Figure. 1 The block diagram of the proposed method

The CPCB data are collected for air quality and other related factors in the monitored station. In this research, the VAR model is applied to perform normalization and augment to make input data suitable for classification. The expanded data were applied to the LSTM model to predict air quality. The main principle of the proposed method is to use the VAR model to perform normalization and augmentation of the input time-series data of pollutants. Then, utilizing it with reduced loss of information to train the LSTM model and better predict AQI Values. The block diagram of the proposed method is shown in Fig. 1.

### 3.1.    Vector Autoregression model

Air quality factors are complex and dynamic relationships are present among the features [16]. The general simultaneous equations model has lower efficiency in revealing dynamic effects for exploringlag phase effects of explained variables of explanatory variables on its own. The available simultaneous equations have set variables as exogenous or endo generous variables that miss someimportant lag variables. Subjective settings in equations model error are reduced by considering allvariables as endogenous in the VAR model [17]. TheVAR model has the following advantages compared to the traditional single equation, (1) generality of the VAR model, and this is easy to add explanatory variablesdue to this is not based on theories. (2) VAR model reveals the short-term and long-termrelationship between air quality factors. The VAR models have limitations, such as many parameters are required to measure, and high correlation in explanatory variables lag periods. The studies show that CO2 emissions and their driving forces have a large number of dynamic relationships. The VAR model is used for dynamic effects analysis of the drivingforce of CO2 emissions.

Eq. (1) provides the formula for the VAR model.

$$y_t = v + A_1 y_{t-1} + \cdots + A_p y_{t-p} + \mu_t \quad ,$$
$$t = 0, \pm 1, \pm 2 \tag{1}$$

where random vector $(K \times 1)$ is in $y_t = (y_{1t}, \ldots, y_{kt})'$, coefficient matrix of $(K \times K)$ is denoted as $A_i$, intercept terms of $(K \times 1)$ vector is denoted as $v = (v_1, \ldots, v_k)'$.

The random error term of K-dimensional is denoted as $\mu_t = (\mu_{1t}, \ldots, \mu_{kt})'$, and classic econometric assumptions are given as

$E(\mu_t, \mu') = 0 \quad (s \neq t)$ and $E(\mu) = 0, E(\mu, \mu') = \sigma^2$.

If there is no further statement, a non-singular matrix is denoted as $\sigma^2$. The resultant vectors of VAR are the augmented and normalized values of pollutant data. These values are utilized as input to LSTM to improve the $\mu$ learning rate and reduce prediction errors.

## 3.2. Long Short Term Memory

The LSTM can retain the important information for the long term based on cell and forget gate. The classification of arrhythmia signals not only requiresrecent data and also previous data. So, the LSTM model has the advantage of handling long-term dependence problems based on a hidden layer of a self-feedback method [16, 17]. Memory cell and

$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t \qquad (5)$$

Three gates, such as input, forget, and output gates, wereused to store information in the LSTM model to help handle the problem of long-term features [18-20].The architecture of the Bi-LSTM model is shown in Fig. 2.
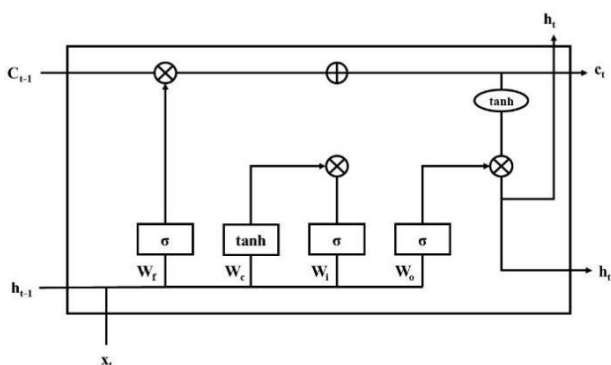


Figure. 2 Architecture of Long Short Term Memory(LSTM) model

The LSTM cell output is denoted as $h_t$, the memory cell value is represented as $c_t$, the previous momentLSTM cell output is denoted as $h_{t-1}$, and the LSTMcell input data is denoted as $x_t$ at time $t$. The LSTM unit calculation process is explained in steps.

1) The candidate memory cell $\tilde{c}_t$ is calculated, the bias is denoted as $b_c$, and the weight matrix isrepresented as $W_c$, as shown in Eq. (2).

$$\tilde{c}_t = \tanh(W_c.[h_{t-1}, x_t] + b_c) \qquad (2)$$

2) The input gate $i_t$ is calculated, current input data update of memory cell state value controls by input gate, the bias is denoted as

$b_i$, the weight matrix is represented as $W_i$, and the sigmoid function is indicated as $\sigma$, as shown in Eq. (3).

$$i_t = \sigma(W_i.[h_{t-1}, x_t] + b_i) \qquad (3)$$

3) The forget gate $f_t$ value is calculated,me forgetting gate controls the memory cell state value based on historical data updng gate, the biasis denoted as $b_f$, the weight matrix is representedas $W_f$, as given in Eq. (4).

$$f_t = \sigma(W_i.[h_{t-1}, x_t] + b_f) \qquad (4)$$

4) The current moment memory cell $c_t$ is calculated, and the last LSTM unit state valueis denoted as $c_{t-1}$, as given in Eq. (5).Where '.' denotes the dot product. Input and forget gate controls update the memory cell based on the state value of the candidate and last cell.

5) The output gate $o_t$ value is calculated, the memory cell state value output is controlled by the output gate, the bias is denoted as $b_0$,

and the weight matrix is represented as $W_0$ as given in Eq. (6)

$$o_t = \sigma(W_0.[h_{t-1}, x_t] + b_0) \qquad (6)$$

6) The LSTM unit output $h_t$ is calculated, as given in Eq. (7).

$$h_t = o_t * \tanh(c_t) \qquad\qquad (7)$$

LSTM model update, reset, read and keep long-time information quickly based on memory cell and control gates. The LSTM model sharing mechanism of internal parameters controls the output dimensionsbased on weight matrix dimensions' settings.

## 3.3. Combination of LSTM and VAR

Neural network training is improved based on thefitted VAR model. Multivariate data of internal behavior of VAR model for adjusting insane values of multivariate data correcting reconstructing NaNs correctly and anomalous trends. Fitness value contains information that is modified original data version manipulated in the model during the training procedure. The kind of augmented data of source original train. The process of the VAR-LSTM model is given in Figure 3.
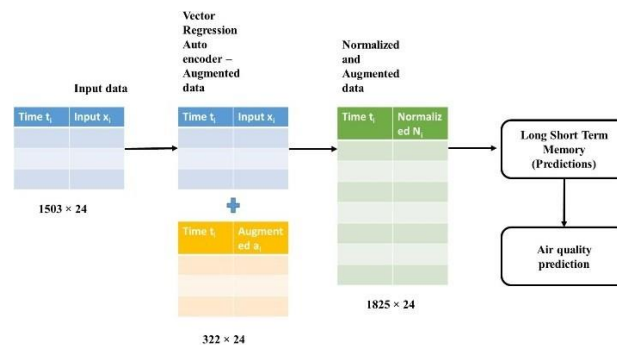


Figure. 3 VAR-LSTM model

A two-step training process is in this strategy. One-step forecasting of all series is carried out usingtfeeding LSTM model and VAR fitted values.Training with raw data is the same differential data tofit VAR. LSTM handles external data sources, for instance, weather conditions or attributes like months,hours, and weekdays to encode cyclically.

A neural network learns from two different data sources and provides better performance on test data.The Vanishing Gradient Problem needs to be handled through multiple steps training. If two tasks are applied to a neural network, the network forgets the first task, and this is a common problem in neural networks.

**Parameter settings:** LSTM model has parametersettings of 0.01 learning rate, 20 epochs, 0.2 dropoutrate, and Adam optimizer is used.

**Metrics:** The formula for Mean Squared Error (MSE), Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and Accuracy is given in equations (8 - 11).

To solve this problem, the entire network is

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y_i})^2 \qquad\qquad (8)$$

It is needed to be appropriately tuned to provide a benefit in performance. The final part of previous training is considered as validation from the observation.

The network is simple, Timeseries Generator of

$$RMSE = \sqrt{MSE} \qquad\qquad (9)$$

Keras is used to fit the model, and Keras-hypetune is

$$\sum^n$$

$$|yi-xi|$$

used for training. Neural Network structures ofhyperparameter optimization are carried out by this

$$MAE =$$

$$i=1$$

$$n \quad (10)$$

framework in a very intuitive way. Some parameter combinations of random search are conducted. All three training involved such as standard fit on raw data, fit on raw data and fit on VAR fitted values.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$

## 4.1. Data Collection

$\times 100 \quad (11)$

The network trained on fitted value of VAR and network trained on original training data are compared. Errors of MAE form are used for output series which are lower for two training steps. The correlation of prediction and actual with 1 delay lag is maintained under 80 % and this is good practice toverify if future predictions are not present values repeated i.e., not a useful prediction.

## 4. Results

The proposed VAR-LSTM method is applied for

the air quality prediction and tested its performance.

The publicly available Indian air quality database wasused to test the performance of the proposed and existing method. The data analysis is performed to understand the features of the dataset. This section provides a detailed description of the dataset analysisand the result of a proposed and existing method.
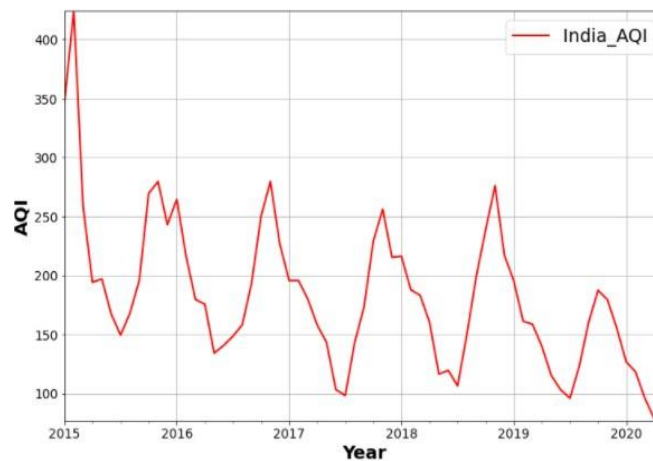
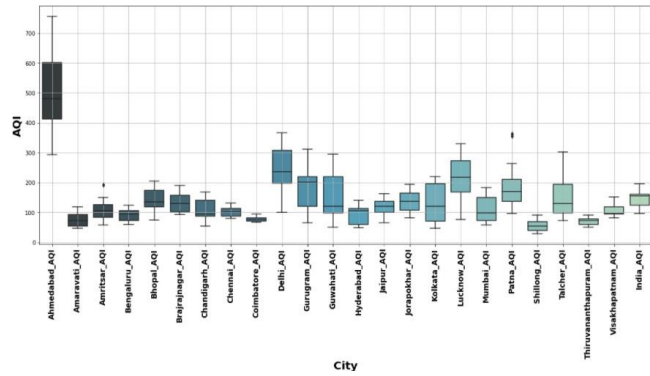

Figure. 4 India AQI for 2015 - 2020 years

Figure. 5 AQI for various cities in India

Fig. 4 provides the details of the India AQI for 2015 – 2020 years, as this shows 2015 has higher AQI and 2020 year has lower AQI. The AQI isslightly increased in 2019 year compared to 2018 year.

Fig. 5 shows the various cities AQI in India, as this shows that Ahmedabad has higher AQI and Delhi has the second higher AQI. The Coimbatore has less AQI in the graph and Shilong city has the second lower AQI in India.

## 4.2. Evaluation

The proposed VAR-LSTM method error value for various epochs values is shown in Fig. 6. The error value is significantly reduced after 50 epochs and maintained the error in the model.
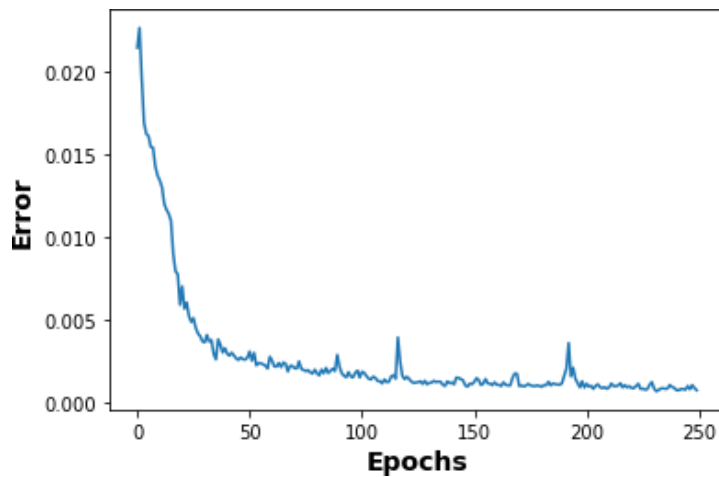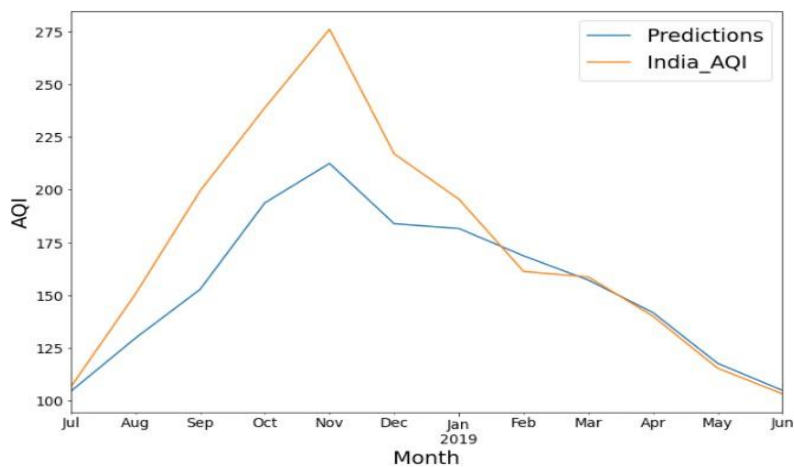


Figure. 6 Error value for various epochs



AQI

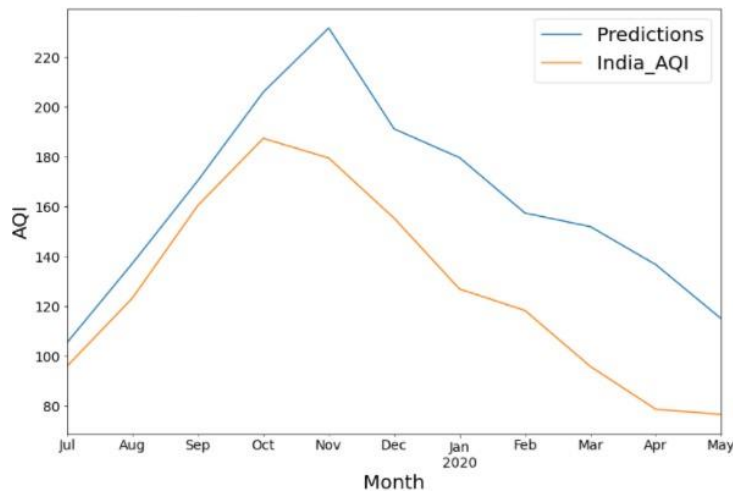Figure. 7 The proposed method predictions and actual India AQI value

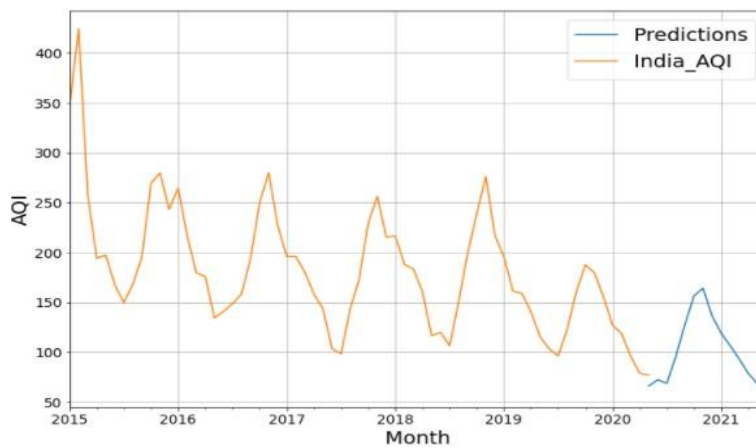Figure. 8 The proposed method and actual AQI in 2020 year



Figure. 9 Training and test data of AQI data

The proposed method predictions in the 2019year in month wise is compared to the actual value,

as shown in Fig. 7. The proposed method has a highererror value from September to January. The proposed

method has higher prediction performance fromMarch to June.

The proposed VAR-LSTM method predictionerror is compared with actual AQI in 2020 year, as

shown in Fig. 8. This shows that the proposed methodhas a higher prediction error from November to May.The proposed method has higher predictionperformance in July to November.
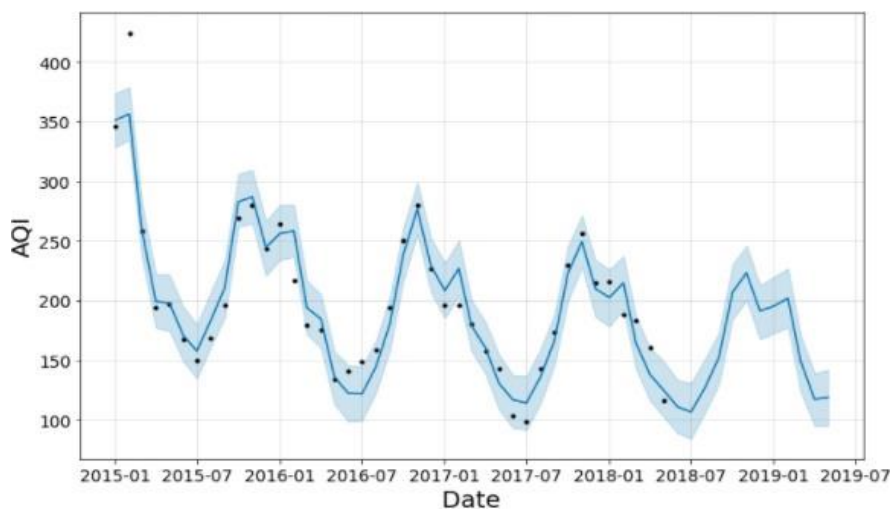


Figure. 10 Prediction performance in time series

Table 1. Performance and error metrics of the proposed method in prediction

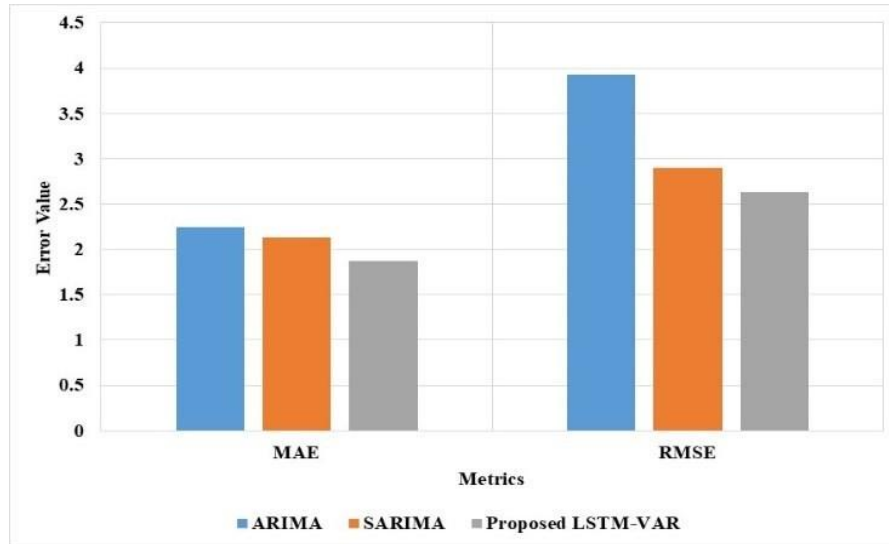|  | MAE | MSE | RMSE | Accuracy |
|---|---|---|---|---|
| ARIMA | 2.245 | 15.46 | 3.932 | 89.31 |
| SARIMA | 2.132 | 8.386 | 2.896 | 91.65 |
| Proposed VAR-LSTM | 1.872 | 6.932 | 2.633 | 93.56 |



Figure. 11 The proposed and existing method AQI prediction

Table 2. Comparative analysis Performance and error metrics of the proposed method in prediction

|  | MAE | RMSE |
|---|---|---|
| Transfer learning LSTM [14] | 5.6878 | 8.8711 |
| NLSTM [16] | 3.456 | 5.579 |
| LSTM-SVR [17] | - | 10.995 |
| Proposed VAR-LSTM | 1.872 | 2.633 |

The proposed VAR-LSTM method prediction and actual data in Air quality prediction are shown in Fig. 9. The data from the year 2015 to 2020 were used as training data and the years of mid-2020 and 2021 were used as testing data. Fig. 9 shows that the model has higher performance in the air quality prediction.

The proposed method prediction is plotted in the graph from the year 2015 to 2019 for every 6 months, as shown in Fig. 10. The dot in the figure represents the error value of the proposed VAR-LSTM model prediction. This shows that the model has a higher error value in the year 2018 to 2019 due to differences in the value of input data.

The proposed VAR-LSTM model accuracy and error metrics are compared with existing methods of SARIMA and ARIMA model, as given in Table 1. This shows that the proposed method has a lower error rate and higher performance than the existing method due to the VAR model being used for training the LSTM model. The VAR method analysis the multivariate in the data and normalize the data to reduce the difference and keep the model suitable for the LSTM training. The hyperparameter optimization of the LSTM model helps to improve the performance of the prediction process. The ARIMA and SARIMA models have higher error rates due to the lower learning in the feature differences. The proposed VAR-LSTM method has 93.56 % accuracy in the prediction and SARIMA has 91.65 % accuracy.

The proposed method shows better performance in terms of RMSE and MAE compared to recent LSTM based methods as shown in Table 2. As shown Transfer learning LSTM has 5.6878 MAE, NLSTM has 3.456 MAE whereas VAR-LSTM possess only 1.872. And it possesses a 2.633 RMSE which is lesser than both NLSTM and LSTM-SVR. Thus it has the more accurate and stable prediction of AQI compared to other LSTM based technologies based on its improved learning rate.

The proposed and existing method error value of AQI prediction is compared in Fig. 11. The proposed VAR-LSTM method has a lower error rate due to the augmentation and normalization of the VAR method. The VAR method normalizes the data to reduce the differences in features and augment the data to make it suitable for LSTM training. The proposed VAR- LSTM method has 1.872

MAE and the existing SARIMA method has a 2.132 MAE value.

## 5. Conclusion

In this research, the VAR-LSTM model is proposed to improve the performance of the air quality prediction. The VAR-LSTM model has the advantage of normalizing the input data based on

multivariate characteristics and augmenting the data to make it suitable to train LSTM. The India Air Quality Index data in Central Pollution Control Boardwere collected to test the proposed and existing models. The proposed method has higher performance in prediction due to the feature selectionof the VAR-LSTM model. The proposed VAR- LSTM model has 6.932 MSE, the SARIMA model has 8.386 MSE and ARIMA model has 15.46 MSE. The future work of the proposed method involves applying the weather data to reduce the error rate in the Air Quality Prediction.

## Conflicts of Interest

The authors declare no conflict of interest.

## Author Contributions

The paper background work, conceptualization, methodology, dataset collection, implementation, result analysis and comparison, preparing and editingdraft, visualization have been done by first author.

The supervision, review of work and project administration, has been done by second author.

## References

[1]    J. Wang, and G. Song, "A deep spatial-temporal ensemble model for air quality prediction", *Neurocomputing*, Vol. 314, pp. 198-206, 2018.

[2]    D. Zhu, C. Cai, T. Yang, and X. Zhou, "A machine learning approach for air quality prediction: Model regularization and optimization", *Big data and cognitive computing*,Vol. 2, No. 1, pp. 5, 2018.

[3]    Y.C. Lin, S.J. Lee, C.S. Ouyang, and C.H., Wu, "Air quality prediction by neuro-fuzzy modelingapproach", *Applied soft computing*, Vol. 86, pp. 105898, 2020.

[4]    D. Schürholz, S. Kubler, and A. Zaslavsky, "Artificial intelligence-enabled context-aware airquality prediction for smart cities", *Journal of Cleaner Production*, Vol. 271, pp. 121941, 2020.

[5]    X. Xu, and M. Yoneda, "Multitask air-quality prediction based on LSTM-autoencoder model",*IEEE transactions on cybernetics,* 2019.

[6]    Y. Zhang, Y. Wang, M. Gao, Q. Ma, J. Zhao, R. Zhang, Q. Wang, and L. Huang, "A predictive data feature exploration-based air qualityprediction approach", *IEEE Access,* Vol. 7, pp. 30732-30743, 2019.

[7]    B. Liu, S. Yan, J. Li, G. Qu, Y. Li, J. Lang, and

R. Gu, "A sequence-to-sequence air quality predictor based on the n-step recurrent prediction", *IEEE Access*, Vol. 7, pp. 43331- 43345, 2019.

[8]    L. Ge, K. Wu, Y. Zeng, F. Chang, Y. Wang, and

S. Li, "Multi-scale spatiotemporal graph convolution network for air quality prediction", *Applied Intelligence*, Vol. 51, No. 6, pp. 3491- 3505, 2021.

[9]    Z. Qi, T. Wang, G. Song, W. Hu, X. Li, and Z. Zhang, "Deep air learning: Interpolation,prediction, and feature analysis of fine-grained air quality", *IEEE Transactions on Knowledge and Data Engineering*, Vol. 30, No. 12, pp. 2285-2297, 2018.

[10]    H. Liu, Q. Li, D. Yu, and Y. Gu, "Air quality index and air pollutant concentration prediction based on machine learning algorithms", *Applied Sciences,* Vol. 9, No. 19, pp. 4069, 2019.

[11]    J. Ma, J.C. Cheng, C. Lin, Y. Tan, and J. Zhang,"Improving air quality prediction accuracy at larger temporal resolutions using deep learning and transfer learning techniques", *Atmospheric Environment*, Vol. 214, pp. 116885, 2019.

[12]    J. Ma, Z. Li, J.C. Cheng, Y. Ding, C. Lin, and Z.Xu, "Air quality prediction at new stations usingspatially transferred bi-directional long short- term memory network", *Science of The TotalEnvironment,* Vol. 705, pp. 135771, 2020.

[13]    X.B. Jin, N.X. Yang, X.Y. Wang, Y.T. Bai, T.L.Su, and J.L. Kong, "Deep hybrid model based onEMD with classification by frequency characteristics for long-term air quality prediction", *Mathematics*, Vol. 8, No, 2, pp. 214,2020.

[14]    M. Krishan, S. Jha, J. Das, A. Singh, M.K. Goyal,and C. Sekar, "Air quality modelling using long short-term memory (LSTM) over NCT-Delhi,India. *Air Quality, Atmosphere & Health*, Vol. 12, No. 8, pp. 899-908, 2019.

[15] M.S.K. Abhilash, A. Thakur, D. Gupta, and B. Sreevidya, "Time series analysis of air pollution in Bengaluru using ARIMA model", *In Ambient Communications and Computer Systems,* pp. 413-426, 2018.

[16] A. Hernandez-Matamoros, H. Fujita, T. Hayashi,and H. Perez-Meana, "Forecasting of COVID19 per regions using ARIMA models and polynomial functions", *Applied soft computing*, Vol. 96, pp. 106610, 2020.

[17] S.N. Singh, and A. Mohapatra, "Repeated wavelet transform based ARIMA model for very short-term wind speed forecasting", *Renewable energy*, Vol. 136, pp. 758-768, 2019.

[18] A. Sherstinsky, "Fundamentals of recurrentneural network (RNN) and long short-term memory (LSTM) network", *Physica D: Nonlinear Phenomena,* Vol. 404, pp. 132306, 2020.

[19] V.K.R. Chimmula, and L. Zhang, "Time series forecasting of COVID-19 transmission in Canadausing LSTM networks", *Chaos, Solitons & Fractals*, Vol. 135, pp. 109864, 2020.

[20] F. Shahid, A. Zameer, and M. Muneeb, "Predictions for COVID-19 with deep learning models of LSTM, GRU and Bi-LSTM", *Chaos, Solitons & Fractals,* Vol. 140, pp. 110212, 2020.