# Integrated Dual LSTM Model-Based Air Quality Prediction

**Rajesh Reddy Muley, Vadlamudi Teja Sai Sri, Kuntamukkala Kiran Kumar, and Kakumanu Manoj Kumar**

**Abstract** Although air quality prediction is a crucial tool for weather forecasting and air quality management, algorithms for making predictions that are based on a single model are prone to overfitting. In order to address the complexity of air quality prediction, a prediction approach based on integrated dual long short-term memory (LSTM) models was developed in this study. The model takes into account the variables that affect air quality such as nearby station data and weather information. Finally, two models are integrated using the eXtreme Gradient Boosting (XGBoosting) tree. The ultimate results of the prediction may be obtained by summing the predicted values of the ideal subtree nodes. The proposed method was tested and examined using five evaluation techniques. The accuracy of the prediction data in our model has significantly increased when compared with other models.

**Keywords** XGBoosting · LSTM · Accuracy · Prediction

## 1 Introduction

The amount of exhaust gas produced by several factories and automobiles continues to climb as industrialisation levels rise, substantially increasing air pollution. People's daily lives are significantly impacted by air quality. Accurate air quality forecasting has emerged as a key strategy for reducing pollution and raising air quality. Data on air quality has caused great worry throughout the world. For predicting air quality, time series data prediction techniques are frequently employed, along with time series prediction models and conventional machine learning techniques. Some

R. R. Muley (✉)
Department of Information Technology, Lakireddy Bali Reddy College of Engineering, Mylavaram, Andhra Pradesh, India
e-mail: rajeshreddi.m@gmail.com

V. T. S. Sri · K. K. Kumar · K. M. Kumar
Lakireddy Bali Reddy College of Engineering, Mylavaram, Andhra Pradesh, India

methods imitate the temporal and geographical dependency of air quality data concurrently. However, commonly used machine learning techniques frequently exhibit considerable performance variability under various conditions. Numerous variables, including temperature, wind speed, and geographical arrangement have an impact on air quality. As a result, it is challenging to produce certain and precise prediction results using the popular single model prediction method. Our thoughts in this paper are based on a strategy that has recently been discussed in the literature: integrating various models to predict air quality. When compared with current models, the integrated model can greatly increase the ability to forecast. But there is still much to learn about how to combine the benefits of several models depending on the features of the data collection.

## 2   Literature Survey

Petr Hájek et al. in [1] genetic algorithms optimise the input variable sets for each forecast of an air pollutant. Based on information gathered by the Pardubice city monitoring station in the Czech Republic, models are developed to predict the specific air quality indices for each air pollutant. The results show that when the root mean squared error is taken into consideration, individual prediction model compositions outperform single forecasts of the common air quality index. As a result, these models can be used to produce air quality index predictions that are more accurate one day in advance.

In order to avoid air pollution in urban areas and improve the quality of life for city dwellers, Kang et al. [2] highlighted the importance of conducting work on city air quality forecasting. Following that, AQI prediction models based on back propagation (BP) neural networks, genetic algorithm optimisation, and genetic simulated annealing algorithm optimisation are established. Comparing and evaluating the prediction outcomes reveal that the BP neural network based on genetic simulated annealing method has a higher accuracy rate, excellent generalisation capacity, and global search ability.

According to Wang et al. [3], who found that air pollution was becoming more severe, the most significant air pollutant, PM2.5 in aerosols, had a negative impact on people's regular output, way of life, and employment, as well as their health. As a result, the forecasting of PM2.5 concentration has taken on significant practical importance. The study selects real-time air quality data that is released, collects historical monitoring data of air environmental contaminants, normalises the data, and then splits the sample data into the two sets in a suitable ratio to form the training dataset and test dataset.

A key component of a smart city is a system for measuring and forecasting air quality, Mahajan et al. [5]. Making a forecast system with great accuracy and a reasonable calculation time is one of the biggest challenges. In this study, we demonstrate that a variety of clustering algorithms may be used to forecast fine particulate matter (PM2.5) concentrations reliably and quickly. We cluster the monitoring

stations depending on their geographic proximity using a grid-based methodology. Data from 557 stations that have been distributed throughout Taiwan's Airbox device network is used in the tests and evaluation. The accuracy and processing time of the various clustering algorithms are compared in a final study.

## 3 Existing System

Commonly used machine learning techniques frequently exhibit considerable performance variability under various conditions. Numerous variables, including temperature, wind speed, and geographical arrangement have an impact on air quality. As a result, it is challenging to produce certain and precise prediction results using the popular single model prediction method.

## 4 Proposed System

In this work, a prediction approach based on integrated dual long short-term memory (LSTM) models was created to handle the complexity of air quality prediction. First, a single-factor prediction model that can independently forecast the value of each component in air quality data is created using sequence to sequence (Seq2Seq) technology. The multi-factor prediction model is then the LSTM model plus the attention mechanism. The model takes into account the air quality parameters such as the data from nearby stations and the weather. The two models are then combined using the eXtreme Gradient Boosting (XGBoosting) tree.

## 5 System Architecture

See Fig. 1.

## 6 Flow Chart

See Fig. 2.

**Fig. 1** System architecture

## 7 Results

Single-factor model is subsequently improved using the ATTENTION layer to create a multi-factor (combination of LSTM, sequence 2 sequence, and attention). In order to combine both models and improve prediction accuracy, features from the multi-model are extracted and retrained using XGBOOST.

The screen below displays information from the air quality dataset, which was used to construct this project.

The first row of the dataset's screen in Fig. 3 shows its column names, while the following rows show its values. As the training features, we used PM values, the target variable, and others.

**Fig. 2** Flow chart

All three models—single-factor LSTM, multi-factor LSTM with attention, and multi-factor integrated with XGBOOST—have been coded by our team. Below are the code and output screens for all the models we coded in the Jupyter notebook. You can see BLUE colour comments in each screen to learn about code (Figs. 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18 and 19).

Fig. 3  Air quality dataset



Fig. 4  Above screen we are loading required Python classes, LSTM, attention and XGBOOST class

**Fig. 5** Above screen we define function to calculate MAPE and to normalise values and then reading and displaying dataset values



**Fig. 6** Above screen we are processing dataset and then removing irrelevant columns and then splitting dataset into train and test and for testing we used 50 values which considers 1 test data per second for next one hour. These test values you can see in blue colour text

**Fig. 7** Above screen we are defining function to calculate RMSE, MAE, and MAPE values and then plot TEST data air quality and predicted air quality graph



**Fig. 8** Above screen we are training single-factor LSTM model and below is the LSTM model summary
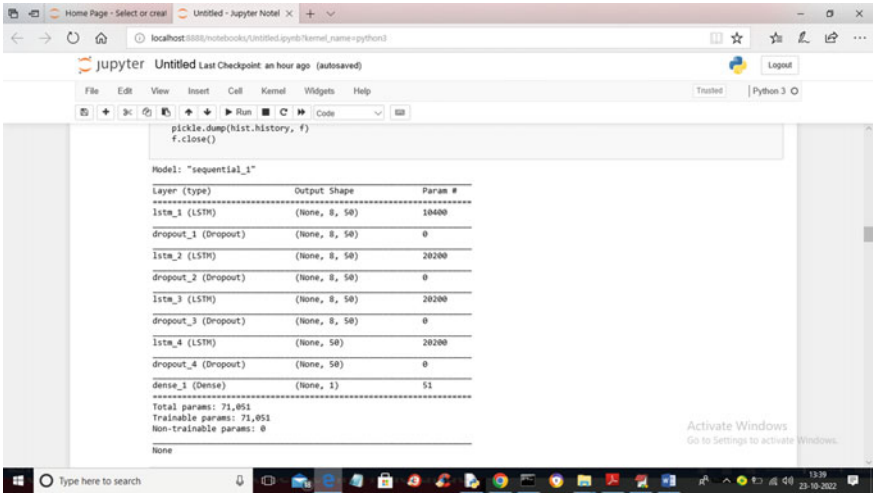
**Fig. 9** Above summary you can see that in single-factor LSTM there is no attention layer and in below screen we can see single-factor model prediction output
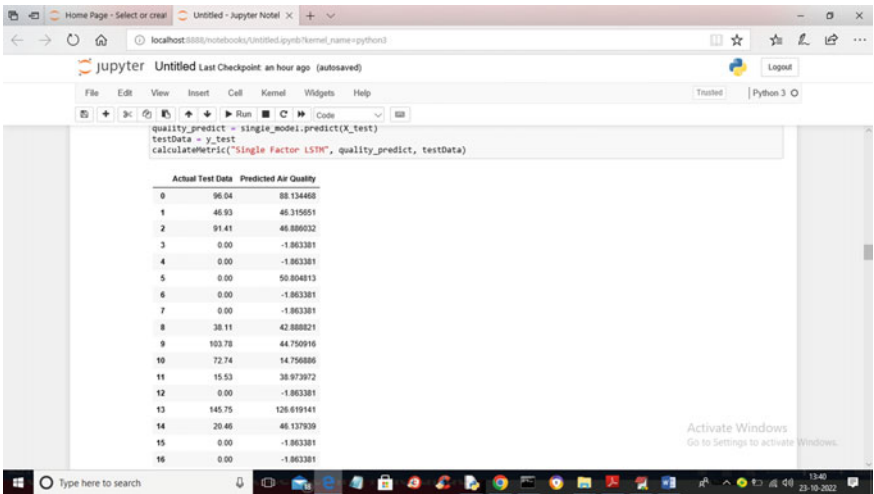


**Fig. 10** Above screen we can see single-factor test data quality and predicted air quality for next 50 records
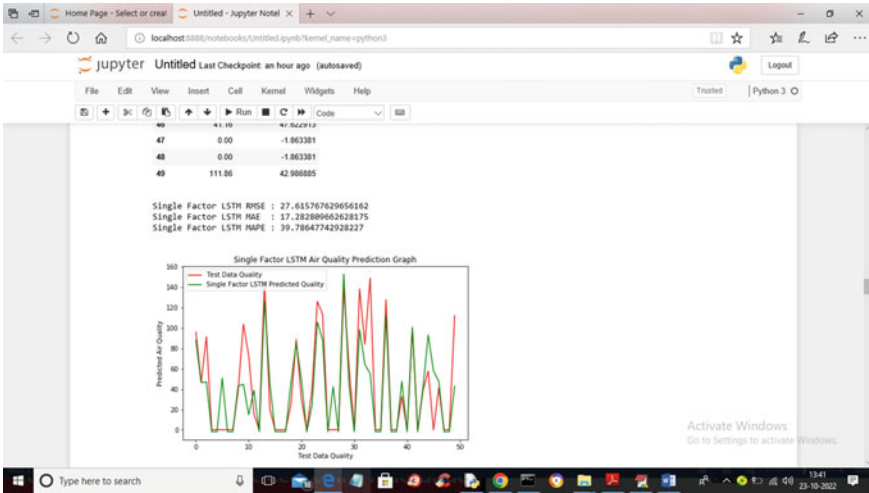
**Fig. 11** Above screen with single factor we got RMSE as 27 and MAE as 17, and MAPE as 39 and then in graph x-axis represents 50 min and y-axis represents air quality values, red line represents TEST data actual values, and green line represents predicted air quality and we can see both lines are overlapping with little difference. In below screen we can see LSTM with attention
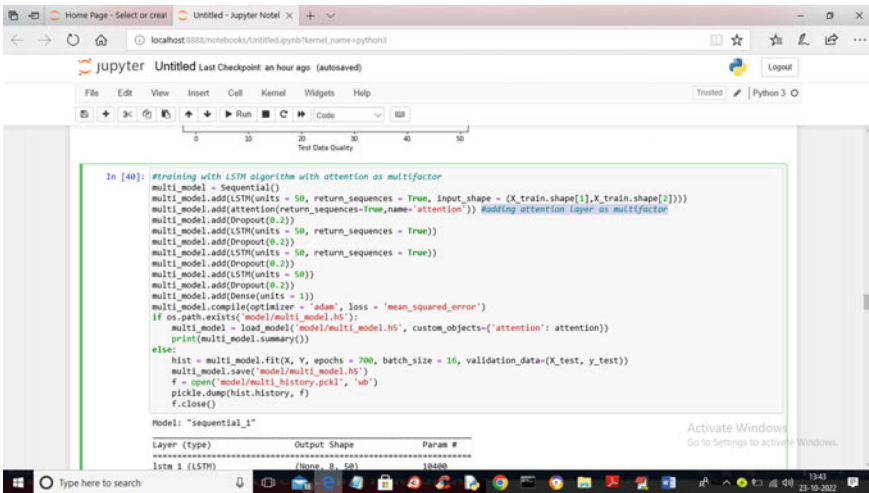


**Fig. 12** Above screen we are defining multi-model by combining LSTM and attention layer and below is the multi-model summary

**Fig. 13** Above screen LSTM is combined with attention and below is the multi-model predicted output



**Fig. 14** Above screen we can see test data and multi-model predicted air quality for next 50 s

**Fig. 15** Above screen with multi-model we got RMSE as 23 and MAE as 13 and MAPE as 37, and we can see both test data and predicted air quality in graph. Above model RMSE, MAE, and MAPE is lesser than single model. In below screen showing integrated model with XGBOOST



**Fig. 16** Above screen showing result of integrated XGBOOST

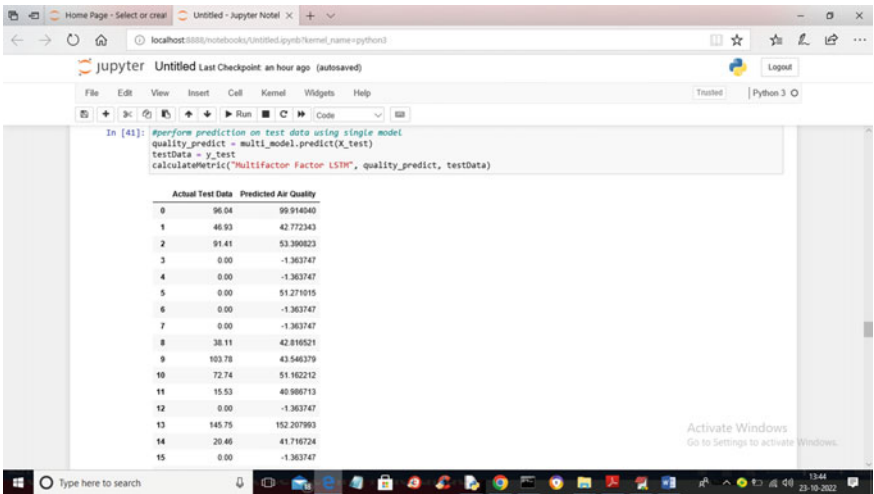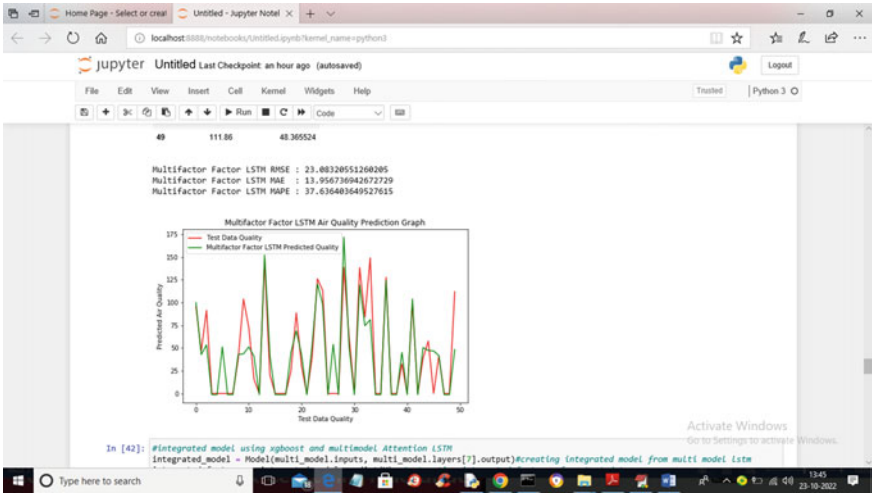**Fig. 17** Above screen with integrated model we got RMSE as 24 and MAE as 12, and MAPE as 31 and we can see predicted and actual test values in graph. In integrated model we got RMSE as high but MAE and MAPE is lesser than single and multi-model. In below screen showing graph of all algorithms



**Fig. 18** Represents values where each different colour bar represents different metric such as RMSE, MAE, and MAPE and in above graph we can see integrated XGBOOST got less MAE and MAPE compared with all other algorithms and same output we can see in below tabular format

**Fig. 19** Above table we can see metric values of all algorithms and integrated XGBOOST got better result with low error rate

## 8  Conclusion

We suggested a prediction model based on integrated dual LSTM model method to increase the precision of air quality data prediction. The integrated model's realisation procedure and impact can be summed up as follows. The air quality characteristics in 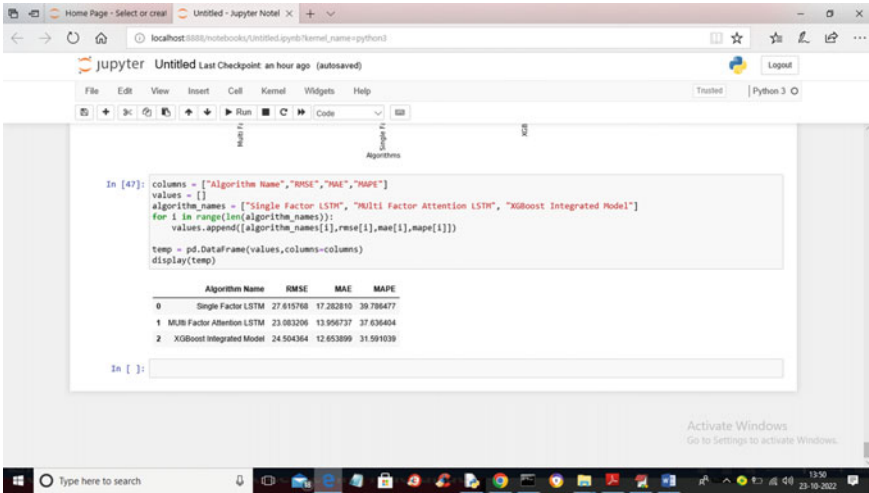the model are taken into consideration, together with meteorological information and data from surrounding stations. The method tree is then used to integrate the two models. First, single-factor models for each factor in the temporal dimension were made. To obtain the forecasted outcomes, the temporal dimension's attributes are employed. The projected value and weight of each leaf node are put together to provide the ideal expected value. Since the technique outlined in this study is based on analysing the experimental data using five evaluation indicators, it can result in predictions that are more accurate.

In order to improve the accuracy of by integrating the advantages of various models, the integrated dual LSTM model technique will be expanded in the next phase of the study. Although our model's outputs have very low probability, we have also found certain prediction results with outlier values. The examination of this sort of outlier value is one of the concerns that has to be addressed in the feature scope.

# References

1. Petr H, Vladimir O (2013) Prediction of air quality indices by neural networks and fuzzy inference systems. Commun Comput Inf Sci 383:302–312. https://doi.org/10.1007/978-3-642-41013-0_31
2. Kang Z, Qu Z (2017) Application of BP neural network optimized by genetic simulated annealing algorithm to prediction of air quality index in Lanzhou. In: Proc. IEEE Comput. Intell. Appl. (ICCIA), Sep. 2017, pp. 155–160. https://doi.org/10.1109/CIAPP.2017.8167199
3. Wang X, Wang B (2019) 'Research on prediction of environmental aerosol and PM2.5 based on artificial neural network.' Neural Comput Appl 31(12):8217–8227. https://doi.org/10.1007/s00521-018-3861-y
4. T. S. Rajput and N. Sharma, "Multivariate regression analysis of air quality index for Hyderabad city: Forecasting model with hourly frequency," Int. J. Appl. Res., vol. 3, no. 8, pp. 443–447, 2017. Accessed: Mar. 20, 2021. [Online]. Available: https://www.allresearchjournal.com/archives/2017/vol3iss ue8/PartG/3–8- 78–443.pdf
5. Mahajan S, Liu H-M, Tsai T-C, Chen L-J (2018) 'Improving the accuracy and efficiency of PM2.5 forecast service using cluster-based hybrid neural network model.' IEEE Access 6:19193–19204. https://doi.org/10.1109/ACCESS.2018.2820164
6. Li R, Dong Y, Zhu Z, Li C, Yang H (2019) 'A dynamic evaluation framework for ambient air pollution monitoring.' Appl Math Model 65:52–71. https://doi.org/10.1016/j.apm.2018.07.052
7. Liu B, Yan S, Li J, Qu G, Li Y, Lang J, Gu R (2019) 'A sequenceto-sequence air quality predictor based on the n-step recurrent prediction.' IEEE Access 7:43331–43345. https://doi.org/10.1109/ACCESS.2019.2908081
8. Gu K, Qiao J, Lin W (2018) 'Recurrent air quality predictor based on meteorology- and pollution-related factors.' IEEE Trans. Ind. Informat. 14(9):3946–3955. https://doi.org/10.1109/TII.2018.2793950
9. Benhaddi M, Ouarzazi J (2021) 'Multivariate time series forecasting with dilated residual convolutional neural networks for urban air quality prediction.' Arabian J. Sci. Eng. 46(4):3423–3442. https://doi.org/10.1007/s13369-020-05109-x
10. Song X, Huang J, Song D (2019) Air quality prediction based on LSTM-Kalman model. In: Proceedings IEEE 8th Joint Int. Inf. Technol. Artif. Intell. Conf. (ITAIC), Chongqing, China, May 2019, 695–699. https://doi.org/10.1109/ITAIC.2019.8785751