

AIR POLLUTION FORECASTING A DEEP LEARNING MODEL BASED ON LSTM TECHNIQUE

Sachin Bhargava, M. Tech. Scholar, Department of Computer Science and Engineering, SIRTE, Bhopal

Prof. Goldy Saini, Assistant Professor, Department of Computer Science and Engineering, SIRTE, Bhopal

Dr. Sneha Soni Associate Professor, Department of Computer Science and Engineering, SIRTE, Bhopal

Abstract:-

During the past few years, severe air-pollution problem has garnered worldwide attention due to its effect on health and well-being of individuals. As a result, the analysis and prediction of air pollution has attracted a good deal of interest among researchers. The research areas include traditional machine learning, neural networks and deep learning. How to effectively and accurately predict air pollution becomes an important issue. In this paper, we propose a deep learning based Long Short Term Memory (LSTM) and gate recurrent unit model. In this new model, we combine local air quality monitoring station, the station in nearby industrial areas, and the stations for external pollution sources. To improve prediction MSE, RMSE, MAE and R square, we aggregate LSTM with GRU models into a predictive model for early predictions based on external sources of pollution and information from nearby industrial air quality stations.

Keywords: -

LSTM, GRU, Air Pollution, Air Quality

I. INTRODUCTION

The environmental surroundings are everything around you - the air, the land, and the rivers and oceans. The atmosphere is composed of several gases, water vapor, and dust particles. The important gases are nitrogen, oxygen-argon, carbon dioxide, neon, helium, hydrogen, ozone, etc. By volume, nitrogen is 78% and oxygen by about 21%. Together, these two gases make 99 % volume of the atmosphere, the rest 1% also important for us. Nitrogen is used by the plants for their survival. But plants cannot take nitrogen directly from the air. Bacteria that live in the soil take nitrogen from the air and change its form so that plants can use it.

Air-Water Cycle: - Air is important for different states of water like ice and water vapor. This water cycle assures the required amount of water for every life.

Air-Carbon Cycle: - Air plays a main role in recycling Carbon Dioxide (CO₂) released in the air by breathing and fossil decomposition. The plant produces energy and releases oxygen from CO₂ by the process of photosynthesis. Humans and animals eat plants for the energy required for living. Decomposition of the body after their lives results in CO₂ emission back into the atmosphere.

Air-Temperature Cycle: - The average earth temperature will be down to freezing point without air.

Air-Life Safety Cycle: - Earth atmosphere will save us from harmful radiation and maintain the temperature conducive for living. Air reduces the risk of space rocks reaching earth as it is vaporized in the air.

Air-Sound Cycle:- If air does not exist we could not hear a screaming jet engine nearer to your ear. Sound can be heard as the air carries sound waves from one point to other.

II. AIR POLLUTION

Currently, air pollution is recognized as a significant public pathological condition, responsible for an increasing range of health impacts which are well documented from the outcomes of extensive studies in several areas of the world. While there's little doubt that fast urbanization means we tend to be currently exposed to unhealthy concentrations and additional numerous ambient air pollutants. X-ray radiation imaging studies on the bodies of ancient mummies have detected proof of respiratory illness, emphysema, respiratory organ lump and arteriosclerosis [1, 2], while autopsies have represented in depth carbon deposits within the respiratory organ. This, in turn, has led to a speculative link to the daily inhalation of smoke in confined areas from fuels used for heat, cooking, and lighting [3]. Air pollution is the release of natural or human made harmful gases and particles into the environment. Air pollution has a higher impact on society and threatens humanity's ability to survive. There was a substantial rise in the use of coal in factories and homes during the urbanization and industrial revolution. These outcomes in smog which caused dismalness and mortality in the stale climatic conditions. During the 1952 Great London smog, the tragedy of losing 4000 life in heavy pollutions highlights the connection between air pollution and human health. In this way air pollution is a developing issue in the urban locales around the globe [4]. Air pollution is a composite of either natural or human activity gasses or particles released into the atmosphere. The number of particles released will be more harmful than the tolerable amount. There are two types of pollutant sources. Natural Sources: Natural pollutants are harmful substances which are emitted by natural phenomenon. SO₂, CO₂, NO₂, CO, and Sulphate are few natural pollutants discharged due to eruptions of volcanoes and forest wildfires.

Implications of Air Pollution to the Quality of Living: Air pollution is a substantial threat to health worldwide. As indicated by 2015 Global Burden of Disease exposure to external pollution is that the world's fifth major death risk problem, accounting although air pollution is a universal problem, it is probable to cause the biggest damage in sensitive people exposed to harmful pollutants. People with chronic diseases (especially cardiorespiratory diseases), very little social support, and lack of medical facilities are most at risk from pollution. In 2015, 4.2 million fatalities and losses of 103.1 million life-years adjusted for disability [8]. The related studies have shown that the related to air pollution for chronic obstructive pulmonary disease and lung cancer varies with age, and these results are biologically possible. The most worrying sign is that the incidence of chronic obstructive pulmonary disease and lung cancer is likely to be higher in older populations (aged >55 years) than younger populations (aged also liable for the depletion of the ozone layer that causes Ultra Violet rays to penetrate the Earth and acid rain that has adverse effects on trees and wildlife [9]. Hence, regulation of air quality and its forecasting has become an important task for both developed and developing nations. As a result of increasing man-made developments, growth of pollutants concentration into the atmosphere has become inevitable and thus, depreciating air quality. Air pollutants are classified into two categories – primary pollutants and secondary pollutants. Pollutants which are generated through the process, for instance, ash from a volcanic eruption are referred to as primary pollutants. Examples – Carbon Monoxide (CO), Sulphur Dioxide (SO₂). Secondary pollutants are a result of the direct or indirect reaction of primary pollutants. A prominent example of secondary pollutant includes ground level Ozone (O₃). Among the six criteria pollutants, Particulate Matter 2.5 (PM_{2.5}) is considered as one of the most pernicious (Pandey et al., 2013). With only 2.5 microns in diameter in size and light advocating them, these tiny particles tend to stay for an extended period of time in the atmosphere and cause harmful effects inside filters of the nose and throat. Growing lung cancer mortality and 4-8 percent increase in the threat of cardiopulmonary is directly associated with the upgradation of each 10- µg/m³ long-term average PM_{2.5} (Pope III et al., 2002). A major source of Nitrogen Dioxide (NO₂) formation is emissions of power plants and automobiles which in turn leads to the formation of ground-level O₃ and fine particle pollution. The mixture of O₃ and NO₂ is considered as a major threat to

children and people suffering from lung diseases like chronic bronchitis, asthma, emphysema [10]. Prolonged exposure to a certain concentration level of O₃ can also result in detrimental effects on plants, crop yield, flora, and fauna. Natural and anthropogenic emission sources of CO include forest fires, animal metabolism, IC engines and burning of carbon enriched fuels. It directly leads to the condition of subnormal oxygenation of the arterial blood and augmentation of greenhouse gases.

III. AIR QUALITY INDEX

An AQI is a measured metric used to record and report evenly on the air quality of various constituents in terms of human health [11]. AQI is a daily air quality reporting index. It informs you how well we are breathing clean air. The AQI is being used as a standard for the quality of air. This AQI gives information to the people about the environment in which they are living. This gives the necessary information about the pollution caused due to the emissions from the industries and can act as a feedback factor so as to modify the emissions from those industries and also helps in allocation of funds to the air pollution control boards. This AQI helps us to rank the different locations in order of the pollution they have, thus highlighting the areas which are more polluted and the frequency of potential hazards. AQI helps to determine the changes in air quality that have happened over a given period of time, allowing for the prediction of air quality and control measures for pollution. This AQI has been calculated by significant levels of air pollutants. The calculation of AQI differs from nation to nation, based on the significant pollutants involved. India AQI, as per Central Pollution Control Board (CPCB) notified (<http://www.cpcb.nic.in>) AQI is constituted by SO₂, O₃, CO, PM (included PM₁₀ and PM_{2.5}), and NO₂, Lead (Pb), Ammonia (NH₃), pollutants.

Table 1: AQI Category

AQI Category	PM ₁₀ 24 Hr µg/m ³	PM ₂₅ 24 Hr µg/m ³	NO ₂ 24 Hr µg/m ³	O ₃ 24 Hr µg/m ³	CO 24 Hr µg/m ³	SO ₂ 24 Hr µg/m ³
Good (0-50)	0-50	0-30	0-40	0-50	0-0.1	0-40
Satisfactory (51-100)	51-100	31-60	41-80	51-100	1.1-2	41-80
Moderate (101-200)	101-250	61-90	81-180	101-168	2.1-10	81-380
Poor (210-300)	251-350	91-120	181-280	169-208	10-17	381-800
Very Poor (310-400)	351-430	121-250	281-400	209-748	17-34	801-1600
Severe (401-500)	430+	250+	400+	748+	34+	1600+

Challenges and Limitations

1. The pollutants types and levels will vary from one location to another. The sources of pollutants are also different like Natural and manmade.
2. Each pollutant will have an adverse effect on human beings.
3. The monitoring stations will give the measurements of various pollutants. More commonly the data available are 24 Hrs average.
4. All the pollutants have to consider for calculating AQI for a particular location. Most of the pollutants will vary with time. Handling a huge amount of data will be a tough task.
5. More research has been done on predicting the individual forecasting of pollutants but not on the AQI.
6. The information is supplied straight without scrutiny from the analyzers for real time AQI, so it may not be for a statutory purpose.
7. Monitoring and subsequent AQI dissemination involves various steps including the operation of sensors and analysers, their calibration, local server data acquisition, transmission via the Internet to a

central database, etc. Due to multiple technical and operational elements such as lengthy power cuts and maintenance issues, monitoring station functioning may also be impacted. Given these constraints, some interruption in the continuous flow and dissemination of information may occur. However, in the event of breakdowns, immediate action is taken to restore the system to operation within a reasonable period of time.

Need for Air Quality Forecasting

Because of restricted resources and practical execution, an alternative approach to tracking air quality is needed to estimate roughly the temporal and spatial distribution of pollutants. Air Quality models are used to indicate air quality standards. It is the least expensive techniques. Regulatory officials can use this modeling as monitoring instruments to evaluate the impact of emissions on ambient air quality. This can also be used to decrease the emissions required to meet the requirements. Air Quality models are generally mathematical descriptions of pollutant transport, diffusion, and chemical reactions from the sources of pollutants [12]. They accommodate one or a lot of mathematical formulae that include parameters that have an effect on concentrations of pollutants at varying distances downwind of emission sources. Typically, they care for sets of pollutant input data that characterize the emissions, meteorology, and topography of a section and turn out outputs that describe that basis's required air quality. Based on the vital input variable treated the models can be simple or advanced. Advanced models are essentially suited for photochemical air pollution, dispersion in complex terrain, and long-range transport of pollutants. Simpler models are suited for the prediction of particulate matter pollutants of downwind sources.

Air Quality prediction has a higher degree of uncertainties than another forecast, as the forecast must diagnose in addition to standard meteorological variables. The forecasting models reduce the uncertainties by "anchoring" the forecast with prior information available and by 'adjusting' the model with additional information like input variables (past measurement values, pollutant concentrations) and meteorological parameters [13]. All the statistical forecasting methods must be calibrated to the effect of large scale emissions. CO, fine particulate matter with an aerodynamic diameter of less than 10 μm (PM10) and less than 2.5 μm (PM2.5) are the most prevalent predicted pollutants.

IV. PROPOSED METHODOLOGY

Machine Learning (ML) is the data handling frameworks, which are built and executed to show the human cerebrum. The main object of the ML analysis is to develop a process of computing device for modeling the brain to execute various process of computing tasks at a faster rate than the traditional systems. ML execute various tasks such pattern matching and classification, optimization function and data clustering. These errands are exceptionally troublesome for conventional PCs, which are quicker in algorithmic procedure of registering undertakings and exact number juggling tasks. ML gangs substantial number of exceedingly interconnected preparing components called hubs or unit or neuron, which more often than not work in parallel and are arranged in standard models [8]. Every neuron is associated with the other by an association interface. Every association interface is related with weights, which contain data about the info flag. This data is required by neuron net to take care of a specific issue. ML, s aggregate conduct is portrayed by their capacity to learn, review and sum up preparing examples or information like that of human mind.

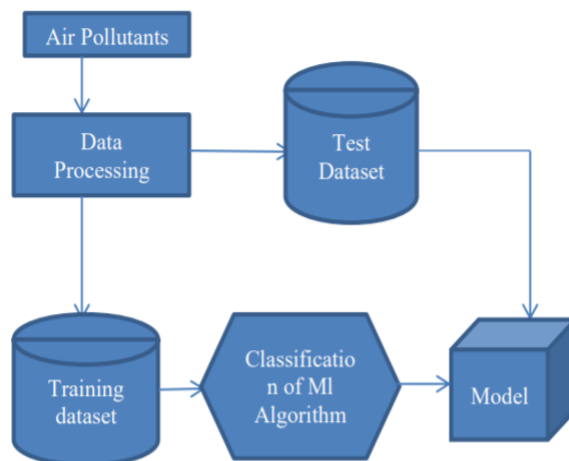


Fig. 1: Basic Diagram of Air Pollution Prediction

Research Methodology

Neural networks normally function like a black box where the decisions are made based on given inputs. It uses static memory in the form of weights to store information about learning experiences. In order to provide explicit representation for memory in RNNs, LSTM network was introduced.

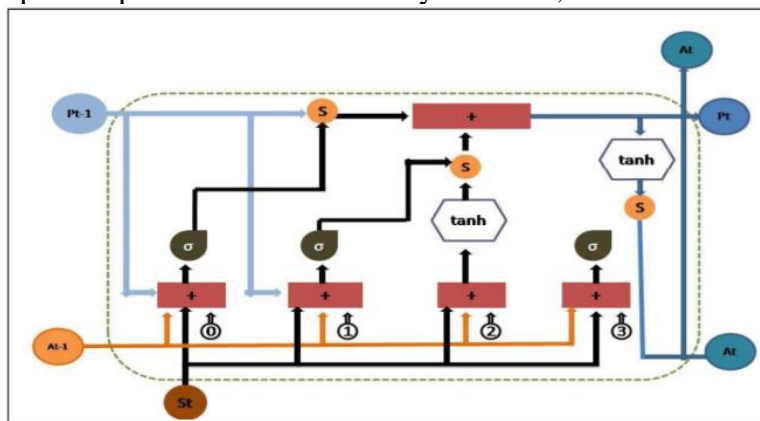


Fig. 2: Working of LSTM

The memory unit is described as 'cell' in the network and these models are adaptation of RNN's and are best suited for sequential data. In this proposed research, we want to investigate the effectiveness of LSTM for sentiment classification of short texts with distributed representation in social media using Evolutionary algorithm. The working of the algorithm is given in fig. 2.

According to above fig. 2 the LSTM network takes three inputs at, ' S_t ', ' A_{t-1} ', ' P_{t-1} '. ' S_t ' is the input vector for the current time step. ' A_{t-1} ' is the output or hidden state transferred by the previous LSTM unit. And ' P_{t-1} ' is the memory element or cell state of the previous unit. It has two outputs such as, ' A_t ' and ' P_t ', where, ' A_t ' is the output of the current unit and ' P_t ' is the memory element of the current unit. Every decision is made after considering current input, previous output and previous memory information. When the current output is obtained the memory is updated. The 'S' indicates the 'Forget' element of multiplication. When the value for the forget element is given as '0' it forgets ninety percent of old memory. For all other values such as 1, 2, and 3 a fraction of old memory is allowed by the unit. The plus operator is present for the piece wise summation to summarize old and new memory. The amount of old memory is decided by the 'S' sign. As a result of two operations, P_{t-1} is changed to P_t . The activation functions described in the fig. 3 are the sigmoid and tanh activation functions having output as a forget valves. The second activation valve is termed as new memory element as it includes old memory while processing new inputs. The old memory, previous output and current input along with a bias vector decides the amount of memory to be given as input to the next unit.

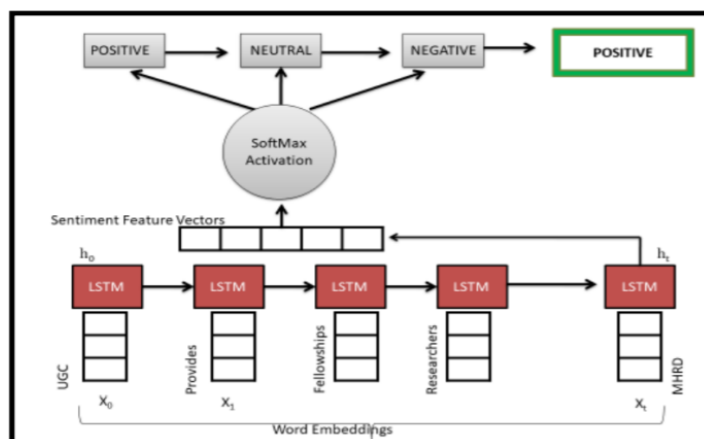


Fig. 3: Sentiment Classification using LSTM

Algorithm: Sentiment analysis on Twitter data using Enhanced LSTM Input: Twitter data set with class labels Output: Classification of tweets whether tweet implies positive, negative or neutral sentiment.

- Step 1: Pre-processed tweets taken in the form of .csv file, data set is loaded
- Step 2: Tagging the tweets
- Step 3: Tagged tweets converted into vectors (word2vector conversion)
- Step 4: Apply Evolutionary algorithm on the vectors to select the best feature set
- Step 5: Enhanced-LSTM performs training only on the best features set selected by Evolutionary Algorithm and obtains a Model
- Step 6: Testing data set is supplied to the Model obtained by Enhanced LSTM
- Step 7: Evaluate the performance of this model based on some parameters

GRU:- Third, the feature sequences is fed into the GRU (GRU) neural networks, which reset gate and update gate constantly adjust their parameters in a large amount of training, so that it can learn the time dependence relationship between the information extracted from the convolutional neural networks. The layer contains only one neuron without any activation function, generating the predicted value of the PM2.5 concentration. Theoretically, the innovation of this method is the combination of the local feature extraction ability and lightness of convnets with the time series prediction ability of GRU by using 1D convnet as a preprocessing step before a GRU. On the other hand, by processing a sequence both way, a bidirectional GRU is able to catch patterns that may have been overlooked by a one-direction GRU.

Table 2: Hyper Parameters

Model	Sequential and RNN
Layers	LSTM, Dense, GRU
Optimizer	Adam
Activation	Relu
Learning rate	0.001
Loss	MSE
Metrics	MSE
Epochs	100

V. SIMULATION RESULT

Results of Air pollution India Dataset

1. Collect data form CPCB website contain 45039 samples with 11 columns namely PM2.5, PM10, NO, No2, Nox, NH3, So2, CO, Ozone the target variable is PM2.5.

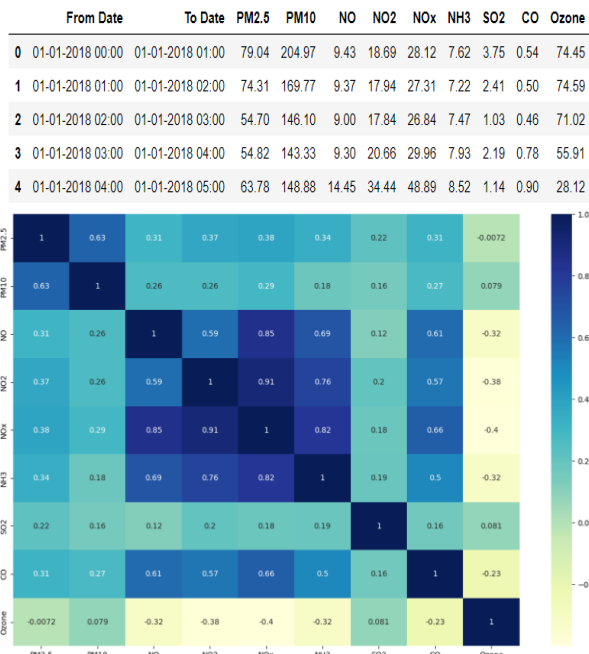


Fig. 4: Exploratory Data Analysis

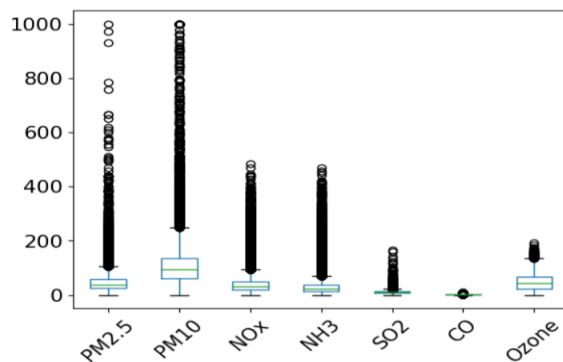


Fig. 5: Different Air Quality Index

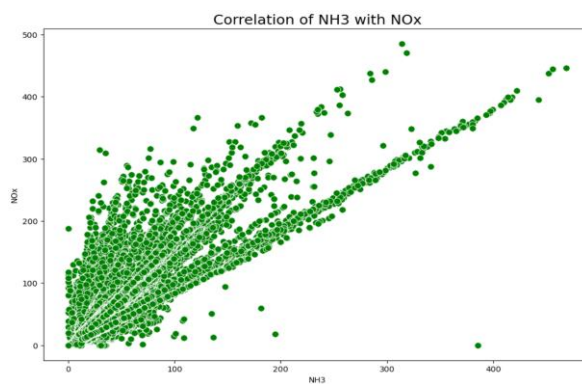


Fig. 6: Correlation of NH3 with Nox

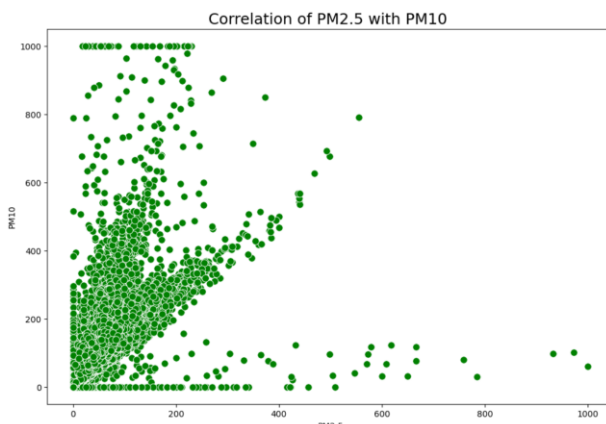


Fig. 7: Correlation of PM2.5 with PM10

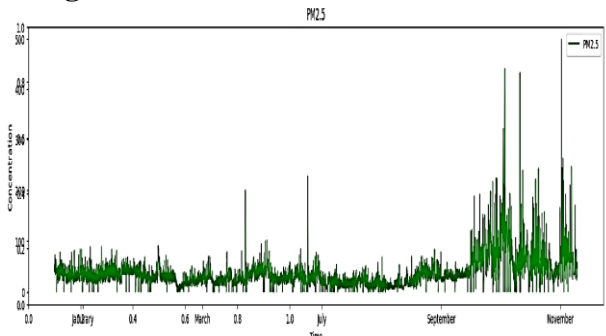


Fig. 8: Concentration of PM2.5 for Month Wise

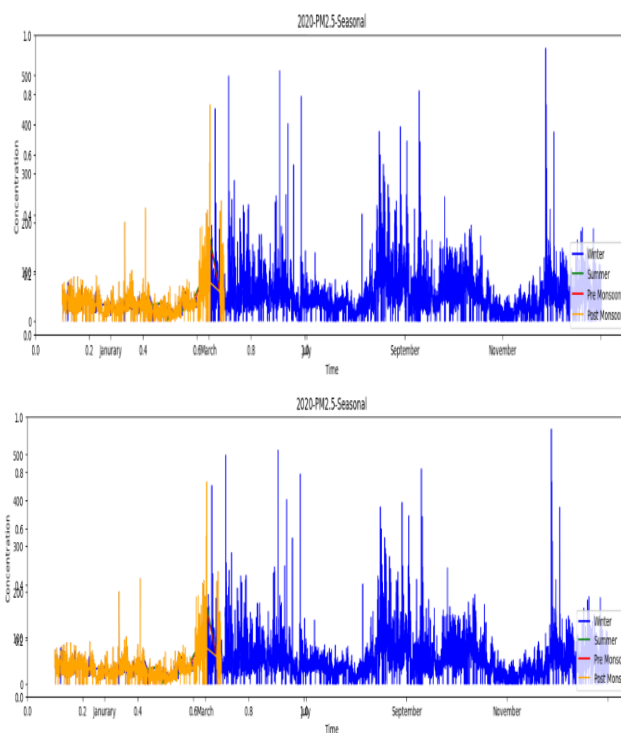


Fig. 9: Concentration of PM2.5 for Seasonal Wise

Table 2: Comparison Result

Models	MSE	RMSE	MAE	R Square
Previous SVM	997.74	31.58	14.35	0.33
Previous LSTM	411.70	20.29	11.64	0.67

Proposed LSTM with GRU	370.62	19.25	11.45	0.70
------------------------	--------	-------	-------	------

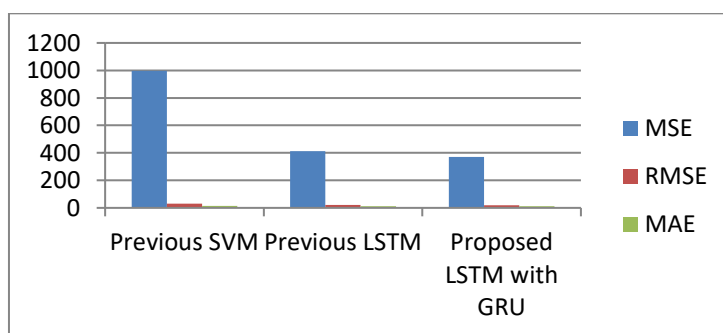


Fig. 10: Graphical Represent of Error

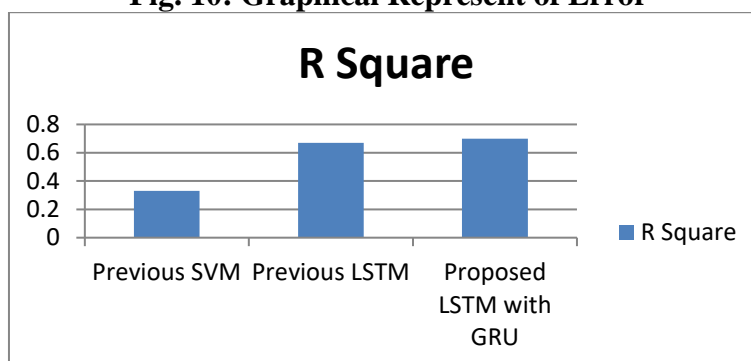


Fig. 11: Graphical Represent of R Square

VI. CONCLUSION

An integrated database of monitoring station air quality and meteorological data can be used for validating the proposed models. The aim of this thesis is to use CBLSTM to investigate a dataset of air pollutants records for the Indian meteorological sector. It is more difficult to determine air quality. This research work will attempt to reduce the risk factor associated with forecasting the Air Quality Index (AQI) of India to a safe human level in order to save a significant amount of meteorological time and resources, as well as to predict whether the air quality is bad or good. In addition, compared with the other benchmark models, the error of the LSTM with GRU model is significantly improved, which shows that the convnets can help the GRU to obtain better prediction performance, because convnets uses its local feature learning ability and subsampling ability to obtain a sequence pattern that is more conducive to GRU processing.

REFERENCES

- [1] Chakradhar Reddy K, Nagarjuna Reddy K, Brahmaji Prasad K and P.Selvi Rajendran, "The Prediction of Quality of the Air Using Supervised Learning", 6th International Conference on Communication and Electronics Systems (ICCES), IEEE 2021.
- [2] Jovan Kalajdjieski, Eftim Zdravevski, Roberto Corizzo, Petre Lameski, Slobodan Kalajdziski, Ivan Miguel Pires, Nuno M. Garcia and Vladimir Trajkovik, "Air Pollution Prediction with Multi-Modal Data and Deep Neural Networks", Remote Sensing, 2020.
- [3] Zhang, Y.; Guo, L.; Wang, Z.; Yu, Y.; Liu, X.; Xu, F. Intelligent Ship Detection in Remote Sensing Images Based on Multi-Layer Convolutional Feature Fusion. Remote Sens. 2020, 12, 3316.
- [4] Bai, L., Wang, J., Ma, X., and Lu, H., 'Air pollution forecasts: An overview', International journal of environmental research and public health 15(4), 780, 2018.
- [5] Balakrishnan, K., Dey, S., Gupta, T., Dhaliwal, R. S., Brauer, M., Cohen, A. J. and Sabde, Y., 'The impact of air pollution on deaths, disease burden, and life expectancy across the states of India: the Global Burden of Disease Study 2017'. The Lancet Planetary Health 3(1), e26-e39, 2019.

- [6] Cai, S., Wang, Y., Zhao, B., Wang, S., Chang, X., and Hao, J., ‘The impact of the “air pollution prevention and control action plan” on PM_{2.5} concentrations in JingJin-Ji region during 2012–2020’, *Science of the Total Environment* 580, 197-209, 2017.
- [7] Cascio, W. E., and Long, T. C., ‘Ambient Air Quality and Cardiovascular Health Translation of Environmental Research for Public Health and Clinical Care’, *North Carolina medical journal* 79(5), 306-312, 2018.
- [8] Huang, M., Zhang, T., Wang, J., and Zhu, L., “A new air quality forecasting model using data mining and artificial neural network”, In *6th IEEE International Conference on Software Engineering and Service Science (ICSESS)*, pp. 259-262, 2015.
- [9] Kang, G. K., Gao, J. Z., Chiao, S., Lu, S., and Xie, G., ‘Air quality prediction: Big data and machine learning approaches’, *International Journal of Environmental Science and Development*, 9(1), 8-16, 2018.
- [10] Kanungo, T., Mount, D. M., Netanyahu, N. S., Piatko, C. D., Silverman, R., and Wu, A. Y., ‘An efficient k-means clustering algorithm: Analysis and implementation’, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, (7), 881-892, 2002.
- [11] Kelly, F. J., & Fussell, J. C., “Air pollution and public health: emerging hazards and improved understanding of risk. *Environmental geochemistry and health*”, 37(4), 631-649, 2015.
- [12] Kemp, A. C., Horton, B. P., Donnelly, J. P., Mann, M. E., Vermeer, M., and Rahmstorf, S., ‘Climate related sea-level variations over the past two millennia’, *Proceedings of the National Academy of Sciences*, 108(27), 11017-11022, 2011.