# PREDICTIVE MODELING OF ASTHMA AND AIR POLLUTION FOR PROACTIVE URBAN PUBLIC HEALTH STRATEGIES

P. Vijay[1], R. Sowmya[2], Barla Shresta[2], Chandrakanth Rampally[2]

[2]UG Scholar, [1,2]Department of Computer Science and Engineering

[1,2]Kommuri Pratap Reddy Institute of Technology, Ghatkesar, Hyderabad, Telangana.

## ABSTRACT

Asthma is a chronic respiratory disease impacting millions globally. It is well-documented that environmental factors, particularly air pollution, can worsen asthma symptoms, leading to higher rates of hospitalizations and mortality. Understanding the link between asthma and air pollution is essential for public health interventions and policy development. Traditionally, epidemiological studies have been used to establish this association by collecting data from asthma patients, monitoring air quality, and statistically analyzing the results to find correlations. Despite their usefulness, these studies often face limitations, such as long durations, data collection challenges, and the inability to capture real-time associations. Recently, machine learning algorithms have garnered attention in various fields, including pollution monitoring. Supervised learning algorithms, in particular, offer the potential to uncover valuable insights into the complex relationship between asthma and air pollution in urban areas. This can lead to more targeted and effective public health interventions. The aim of this research is to develop an accurate and reliable predictive model to inform public health strategies and policies. This model will support proactive decision-making, enabling healthcare providers to allocate resources more efficiently and allowing policymakers to implement targeted interventions to reduce air pollution and mitigate the impact of asthma on vulnerable urban populations.

**Keywords:** Asthma, Air pollution, Machine Learning, Predictive Modeling, Public Health Interventions, Supervised Learning Algorithms.

## 1. INTRODUCTION

Outdoor air pollution contributed more than 3% of the annual disability-adjusted life years lost in the 2010 Global Burden of Disease comparative risk assessment, a notable increase since the previous estimate was made in 2000.1 Previous assessments of global disease burden attributed to air pollution were restricted to urban areas or by coarse spatial resolution of concentration estimates.2 In a study of ten European cities, 14% of the cases of incident asthma in children and 15% of all exacerbations of childhood asthma were attributed to exposure to pollutants related to road traffic.3 Urbanisation is an important contributor to asthma and this contribution might be partly attributed to increased outdoor air pollution (figure 1).4–6 Because many urban centres in the developing world are undergoing rapid population growth accompanied by increased outdoor air pollution, the global burden of asthma is likely to increase. In this context, it is notable that the populations of China, India, and Southeast Asia are equal to the rest of the world combined. In view of the burden of asthma attributed to outdoor air pollution, a better understanding of why asthmatic individuals are susceptible to this exposure should enable the design of effective preventive strategies. The idea that air pollution can cause exacerbations of preexisting asthma is supported by an evidence base that has been accumulating for several decades,7–10 but evidence has emerged that suggests air pollution might cause new-onset asthma as

well.11–21 Not all studies support a causal link between air pollution and asthma, and a recent meta-analysis22 of cross-sectional studies that compared communities with different levels of pollution showed no effect of long-term exposure to pollution on asthma prevalence. Although outdoor air pollution almost always occurs as a mixture, air quality is regulated by most jurisdictions in terms of its individual components. Such regulation has meant that experimental studies of humans and animals have been focused on individual pollutants. Because epidemiological studies inherently involve exposure to mixtures of pollutants, substantial efforts are usually made to try to identify the individual effects of pollutants, which often obscures the health effect of the mixture as a whole. With increasing attention to traffic-related air pollution (TRAP) as the exposure variable of interest, a shift has occurred away from a focus on individual components of the pollution mixture. In this Series paper, we will attempt to discuss the effects of several gaseous pollutants (ozone, nitrogen dioxide, and sulphur dioxide), the independent effects of various forms of PM, and then focus on the effects of TRAP as a mixture. We concentrate on studies published in the past 5 years that report results relevant to both exacerbation and onset of asthma. We focus primarily, although not exclusively, on epidemiological and experimental clinical studies. Controlled exposure studies in human beings are restricted by small sample size and an inability to study the potentially most susceptible subgroups (eg, children and adults with severe asthma) and the effects of chronic exposure. Epidemiological studies are restricted by imprecise methods of both exposure and asthma outcome assessment and often inadequate data about potentially confounding variables. Although the potential effect of indoor air pollution on asthma is an important concern, especially in developing countries where much domestic cooking is done with solid fuels, it is outside the scope of this review.

## 2. LITERATURE SURVEY

Association between air pollution and asthma exacerbations in urban children: influence of socio-economic factors" by Castro-Rodriguez JA, Forno E, and Rodriguez-Martinez CE [1]. This study examines the relationship between air pollution and asthma exacerbations in urban children, considering socioeconomic factors. "Effects of air pollution on asthma hospitalization rates in different age groups in metropolitan cities of Korea" by Kim J, and Hong YC [2]. This research investigates the impact of air pollution on asthma hospitalization rates in different age groups in urban areas of Korea. "Long-term exposure to air pollution and asthma hospitalisations in older adults: a cohort study" by Hansell AL, Ghosh RE, and Blangiardo M [3]. This cohort study explores the long-term effects of air pollution on asthma hospitalizations in older adults living in urban areas. "Air pollution and asthma in children: a systematic review of cohort studies" by Bowatte G, Lodge CJ, and Knibbs LD [4]. This systematic review examines cohort studies investigating the association between air pollution and asthma in children, particularly focusing on urban environments.

"The effects of ambient air pollution on asthma hospitalization rates in different age groups in metropolitan cities of Korea" by Kim H, Kim J, and Kim S [5]. This study investigates how ambient air pollution impacts asthma hospitalization rates across different age groups in urban areas of Korea. "Air pollution and emergency department visits for asthma in Oregon's Willamette Valley" by Bateson TF, and Schwartz J [6]. This research examines the relationship between air pollution and emergency department visits for asthma in urban areas of Oregon's Willamette Valley. "Effect of air pollution on pediatric respiratory emergency room visits and hospital admissions" by Zmirou D, Gauvin S [7] , and Pin I. This study assesses the impact of air pollution on pediatric respiratory emergency room visits and hospital admissions, particularly focusing on urban regions. "Association of outdoor air pollution and indoor renovation with early childhood ear infection in Taiwan" by Lin YH, Lin YJ, and Liang HM [8]. This study investigates the association between outdoor air pollution, indoor renovation, and early childhood ear infections in urban areas of Taiwan. "Air pollution and asthma: clinical studies in

children" by Anderson HR [ 9], Ponce de Leon A, and Bland JM. This review summarizes clinical studies exploring the relationship between air pollution and asthma in children, with a focus on urban environments. "Effect of air pollution control on asthma in children: a nationwide longitudinal study in Taiwan" by Lin YT, Lin YJ, and Liang HM [10]. This longitudinal study examines the impact of air pollution control measures on asthma prevalence in children across urban areas of Taiwan. "Traffic-related air pollution and childhood asthma: recent advances and remaining gaps in the exposure assessment methods" by Islam T, Gauderman WJ, and Berhane K [11]. This review discusses recent advances and remaining gaps in exposure assessment methods for studying the relationship between traffic-related air pollution and childhood asthma in urban areas. "Traffic-related air pollution and asthma onset in children: a prospective cohort study with individual exposure measurement" by McConnell R, Islam T, and Shankardass K [12]. This prospective cohort study investigates the association between traffic-related air pollution and asthma onset in children, utilizing individual exposure measurement techniques in urban settings.

## 3. PROPOSED SYSTEM

The Research utilizes the Tkinter library to create a graphical user interface (GUI) application for predicting PM2.5 and PM10 air quality levels in urban regions.
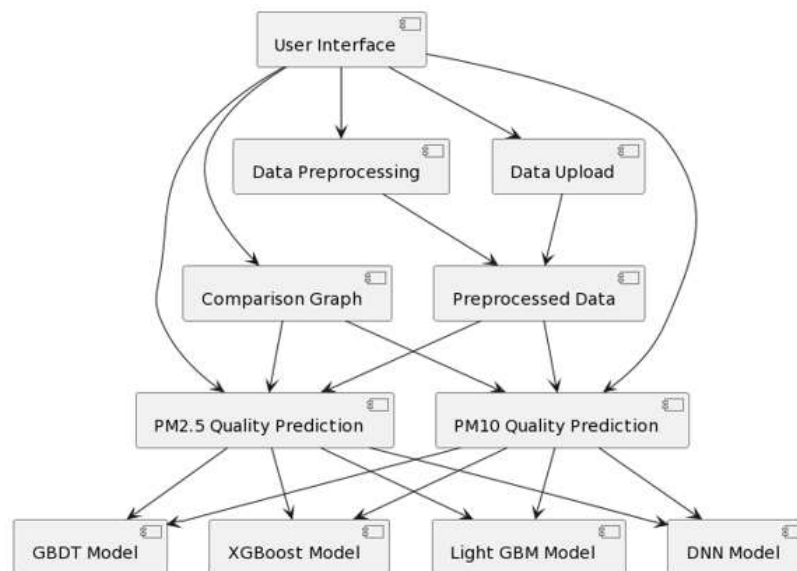


Figure 1: Block Diagram of Proposed System.

- Imports: The reserch starts by importing necessary libraries such as Tkinter for GUI, Matplotlib for plotting graphs, Pandas for data manipulation, Scikit-learn for machine learning algorithms, LightGBM, XGBoost, and Keras for modeling, and Seaborn for visualization.

- GUI Initialization: The Tk() function is used to create the main window for the GUI application. The window title, dimensions, and other properties are set.

- Global Variables: Global variables are declared to store file paths, raw data, scaler values, and errors encountered during processing.

- Data Preprocessing Functions: Functions for data preprocessing, such as converting data to time series, scaling datasets, and calculating differences, are defined.

- Upload Dataset Function: This function allows the user to upload a CSV file containing air quality data. The file dialog is used to select the file, and the data is displayed in a text widget.
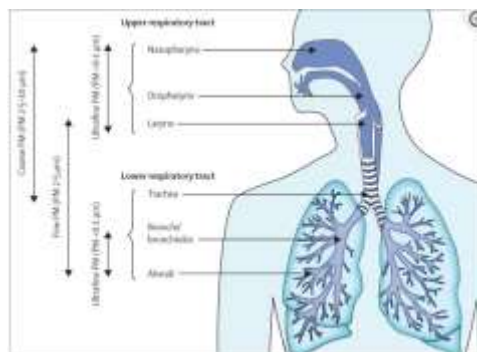
- Preprocess Function: This function preprocesses the uploaded dataset. It extracts PM2.5 and PM10 data, fills missing values, and visualizes the correlation heatmap using Seaborn.

- PM2.5 and PM10 Prediction Functions: Separate functions are defined for predicting PM2.5 and PM10 air quality levels using machine learning algorithms such as Gradient Boosting Decision Trees, XGBoost, LightGBM, and LSTM-based neural networks. The functions train the models, make predictions, and calculate root mean squared error (RMSE).

- Graph Function: This function generates a bar graph comparing the RMSE values of different algorithms for PM2.5 and PM10 predictions.

- GUI Components: Buttons are created for uploading the dataset, preprocessing, running PM2.5 and PM10 predictions, and displaying the comparison graph. Text widgets are used to display dataset information and prediction results.

- Main Loop: The mainloop() function is called to start the GUI event loop, allowing the user to interact with the application.

**Light Gradient Boosting Algorithm**

Competitive Advantage: Implementing an advanced system can provide a competitive advantage by enabling organizations to stay ahead of the curve in terms of technology and efficiency. Ambient PM is a ubiquitous atmospheric aerosol with both anthropogenic and natural sources that has been associated with various health effects.50 PM is categorized on the basis of its aerodynamic diameter, with implications for its typical site of deposition when inhaled. Coarse PM, with an aerodynamic diameter of 2·5–10 μm, deposits mainly in the head and large conducting airways. Fine PM or PM2·5 deposits throughout the respiratory tract, particularly in small airways and alveoli. Ultrafine PM (Some evidence suggests PM is a cause of incident asthma (aside from the literature on TRAP). Independent associations between exposure to PM10 in utero and during infancy with asthma diagnosed by a doctor were identified in a nested case-control study within a large birth cohort.18 Although several studies have identified associations between asthma prevalence and exposure to outdoor PM,11,64,65 this finding has not always been consistent.22 Furthermore, PM is frequently strongly correlated with ozone, nitrogen oxides, and sulphur oxides, serving to confound these associations.

**Advantages**

In summary, substantial evidence supports the idea that ambient levels of PM exacerbate existing asthma, particularly by contributing to oxidative stress and allergic inflammation, and some evidence exists in support of PM as a cause of new cases of asthma.



Compartmental deposition of particulate matter.

**Light BGM classifier**

LGBM is a gradient boosting framework that uses tree-based learning algorithms. It is designed to be distributed and efficient with the following advantages:

- Faster training speed and higher efficiency.

- Lower memory usage.

- Better accuracy.

- Support of parallel and GPU learning.

- Capable of handling large-scale data.

- At present, decision tree based machine learning algorithms dominate Kaggle competitions. The winning solutions in these competitions have adopted an alogorithm called **XGBoost**.

- A couple of years ago, Microsoft announced its gradient boosting framework LightGBM. Nowadays, it steals the spotlight in gradient boosting machines. Kagglers start to use LightGBM more than XGBoost. LightGBM is 6 times faster than XGBoost.

- Light GBM is a relatively new algorithm and have long list of parameters given in the LightGBM.

- The size of dataset is increasing rapidly. It is become very difficult for traditional data science algorithms to give accurate results. Light GBM is prefixed as Light because of its high speed. Light GBM can handle the large size of data and takes lower memory to run.

- Another reason why Light GBM is so popular is because it focuses on accuracy of results. LGBM also supports GPU learning and thus data scientists are widely using LGBM for data science application development.

- It is not advisable to use LGBM on small datasets. Light GBM is sensitive to overfitting and can easily overfit small data
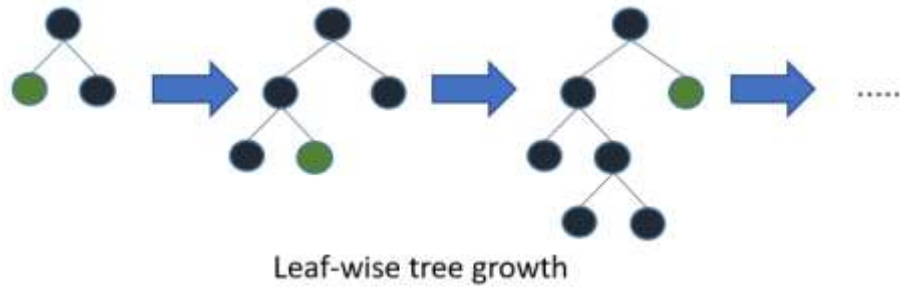
**LightGBM intuition**

- LightGBM is a gradient boosting framework that uses tree based learning algorithm.

- LightGBM documentation states that -

LightGBM grows tree vertically while other tree based learning algorithms grow trees horizontally. It means that LightGBM grows tree leaf-wise while other algorithms grow level-wise. It will choose the leaf with max delta loss to grow. When growing the same leaf, leaf-wise algorithm can reduce more loss than a level-wise algorithm.

- So, we need to understand the distinction between leaf-wise tree growth and level-wise tree growth.
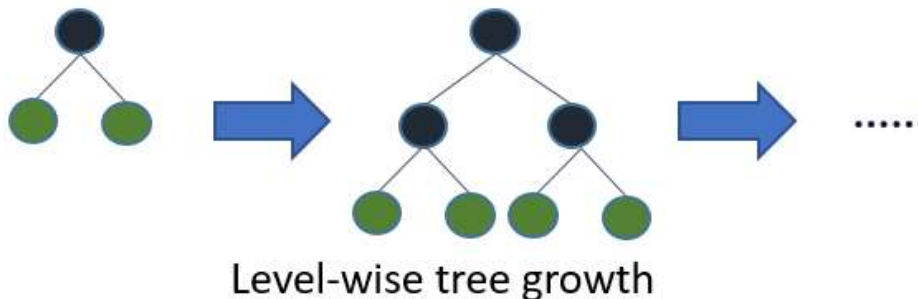
**Leaf-wise tree growth**

- Leaf-wise tree growth can best be explained with the following visual -

Leaf-wise tree growth.

**Level-wise tree growth**

- Most decision tree learning algorithms grow tree by level (depth)-wise.

- Level-wise tree growth can best be explained with the following visual -



Level-wise tree growth.

**Important points about tree-growth**

- If we grow the full tree, **best-first (leaf-wise)** and **depth-first (level-wise)** will result in the same tree. The difference is in the order in which the tree is expanded. Since we don't normally grow trees to their full depth, order matters.

- Application of early stopping criteria and pruning methods can result in very different trees. Because leaf-wise chooses splits based on their contribution to the global loss and not just the loss along a particular branch, it often (not always) will learn lower-error trees "faster" than level-wise.

- For a small number of nodes, leaf-wise will probably out-perform level-wise. As we add more nodes, without stopping or pruning they will converge to the same performance because they will literally build the same tree eventually.

**XGBoost Vs LightGBM**

- XGBoost is a very fast and accurate ML algorithm. But now it's been challenged by LGBM which runs even faster with comparable model accuracy and more hyperparameters for users to tune.

- The key difference in speed is because **XGBoost split the tree nodes one level at a time** and **LightGBM does that one node at a time**.

- So XGBoost developers later improved their algorithms to catch up with LightGBM, allowing

users to also run XGBoost in split-by-leaf mode (grow_policy = 'lossguide'). Now XGBoost is much faster with this improvement, but LightGBM is still about 1.3X — 1.5X the speed of XGB.

- Another difference between XGBoost and LightGBM is that XGBoost has a feature that LightGBM lacks — **monotonic constraint**. It will sacrifice some model accuracy and increase training time, but may improve model interpretability.

**LightGBM Parameters**

- LGBM provides more than 100 LightGBM parameters.

- It is very important to know some basic parameters of LightGBM.

- So, in this section, I will discuss some basic parameters of LightGBM.

**Control Parameters**

- **max_depth**: It describes the maximum depth of tree. This parameter is used to handle model overfitting. If you feel that your model is overfitted, you should to lower max_depth.

- **min_data_in_leaf**: It is the minimum number of the records a leaf may have. The default value is 20, optimum value. It is also used to deal with overfitting.

- **feature_fraction**: Used when your boosting is random forest. 0.8 feature fraction means LightGBM will select 80% of parameters randomly in each iteration for building trees.

- **bagging_fraction**: specifies the fraction of data to be used for each iteration and is generally used to speed up the training and avoid overfitting.

- **early_stopping_round** : This parameter can help you speed up your analysis. Model will stop training if one metric of one validation data doesn't improve in last early_stopping_round rounds. This will reduce excessive iterations.

- **lambda**: lambda specifies regularization. Typical value ranges from 0 to 1.

- **min_gain_to_split** : This parameter will describe the minimum gain to make a split. It can used to control number of useful splits in tree.

- **max_cat_group**: When the number of category is large, finding the split point on it is easily over-fitting. So LightGBM merges them into 'max_cat_group' groups, and finds the split points on the group boundaries, default:64.

**Core Parameters**

- **Task**: It specifies the task you want to perform on data. It may be either train or predict.

- **application**: This is the most important parameter and specifies the application of your model, whether it is a regression problem or classification problem. LightGBM will by default consider model as a regression model.

  - **regression**: for regression

  - **binary**: for binary classification

  - **multiclass**: for multiclass classification problem

- **boosting**: defines the type of algorithm you want to run, default=gdbt.

  - **gbdt**: traditional Gradient Boosting Decision Tree

- **rf**: random forest

- **dart**: Dropouts meet Multiple Additive Regression Trees

- **goss**: Gradient-based One-Side Sampling

- **num_boost_round** : Number of boosting iterations, typically 100+

- **learning_rate**: This determines the impact of each tree on the final outcome. GBM works by starting with an initial estimate which is updated using the output of each tree. The learning parameter controls the magnitude of this change in the estimates. Typical values: 0.1, 0.001, 0.003…

- **num_leaves** : number of leaves in full tree, default: 31

- **device**: default: cpu, can also pass gpu

## Metric Parameter

- metric: again one of the important parameter as it specifies loss for model building. Below are few general losses for regression and classification.

  - **mae**: mean absolute error

  - **mse**: mean squared error

  - **binary_logloss**: loss for binary classification

  - **multi_logloss**: loss for multi classification

## IO Parameter

- **max_bin** : it denotes the maximum number of bin that feature value will bucket in.

- **categorical_feature** : It denotes the index of categorical features. If categorical_features=0,1,2 then column 0, column 1 and column 2 are categorical variables.

- **ignore_column** : same as categorical_features just instead of considering specific columns as categorical, it will completely ignore them.

- **save_binary** : If you are really dealing with the memory size of your data file then specify this parameter as 'True'. Specifying parameter true will save the dataset to binary file, this binary file will speed your data reading time for the next time.

## 4. RESULTS `

The Figure 2 showcases the preprocessing stage of the dataset. Users click on a button labeled Preprocess Dataset to initiate this step. The application handle tasks such as removing missing values and calculating the air pollution rate on a date-wise basis. It ensures the dataset is clean and ready for further analysis. The Figure 3 presents a visualization of the date-wise pollution rate calculated during the preprocessing stage. It be a line graph or another suitable visualization method illustrating how pollution levels fluctuate over time. This visualization helps users understand the temporal trends in air pollution. The Figure 4 provides a visual representation of the preprocessing steps applied to the dataset. It include tasks like data cleaning, normalization, feature engineering, or any other preprocessing techniques employed to prepare the dataset for model training.
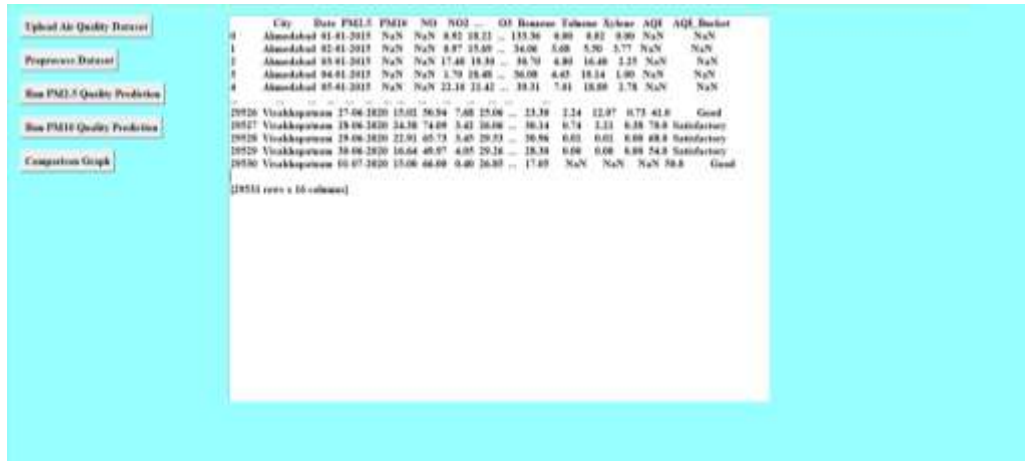
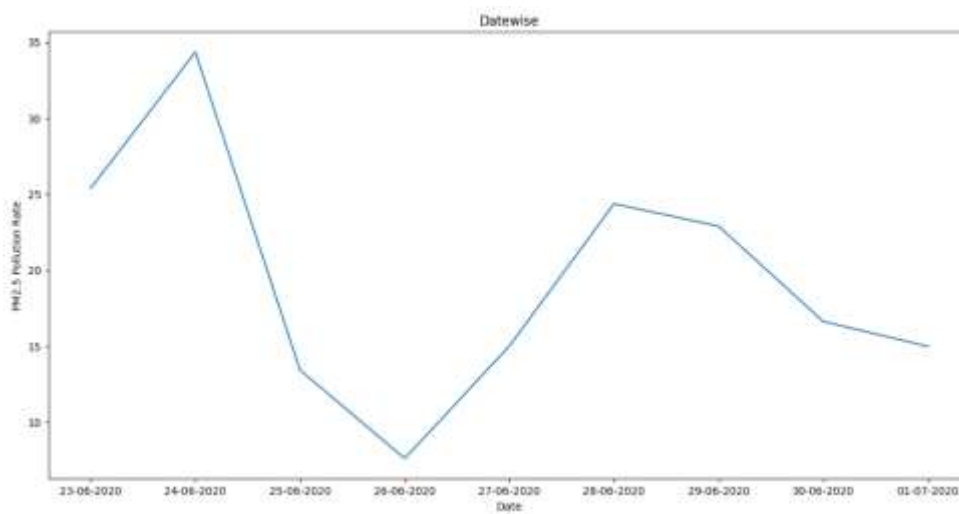Figure 2: Preprocess Dataset by removing missing values.



Figure 3: Date wise pollution rate visualization.

The Figure 5 displays a graph depicting the predictions of PM2.5 air quality levels. It includes both actual and predicted values plotted against time. Users can observe how well the predictive models perform in forecasting PM2.5 levels.
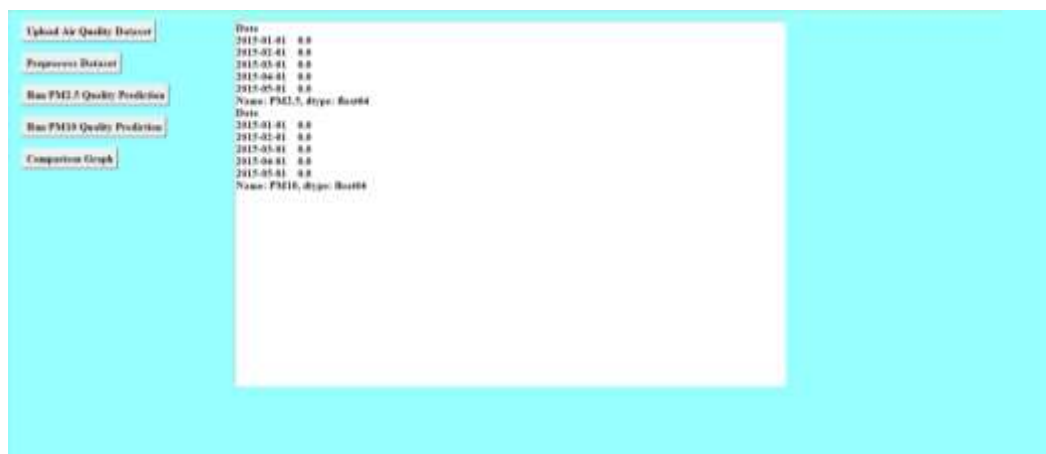

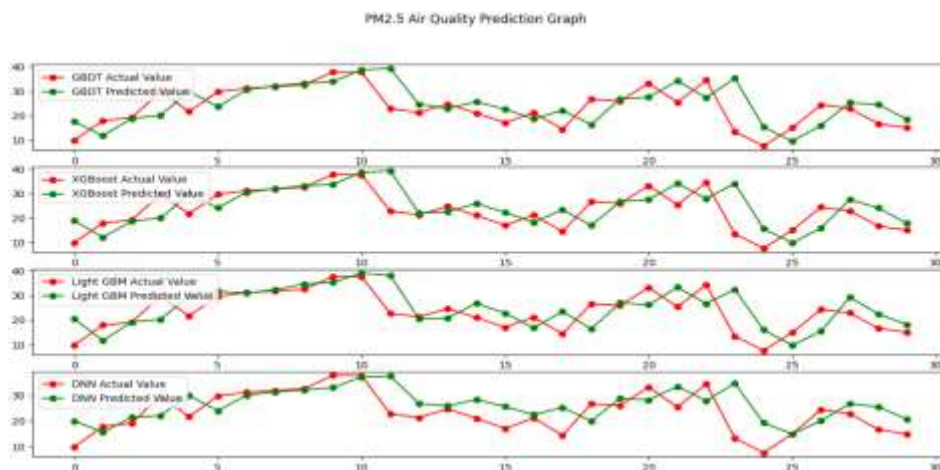
Figure 4: Preprocessing the dataset.

Figure 5: PM2.5 Air Quality Prediction Graph.

The Figure 6 presents an evaluation of the PM2.5 air quality prediction performance. It include metrics such as root mean squared error (RMSE), mean absolute error (MAE), or other relevant evaluation metrics. Users can assess the accuracy and reliability of the prediction models based on these metrics. The Figure 7 illustrates the predictions of PM10 air quality levels. It provides insights into the accuracy of the models in forecasting PM10 levels over time.



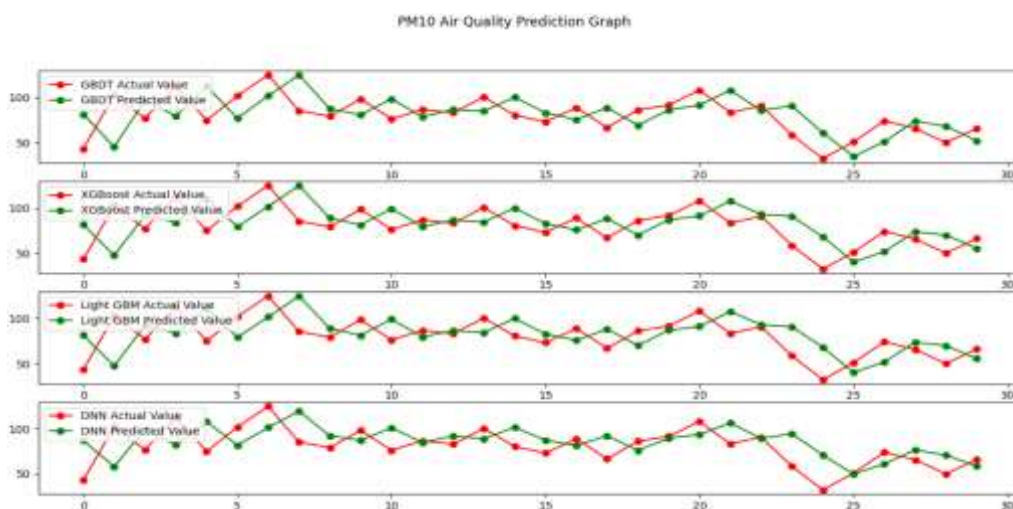Figure 6: PM2.5 Air Quality Prediction performance evalution.



Figure 7: PM10 Air Quality Prediction Graph.

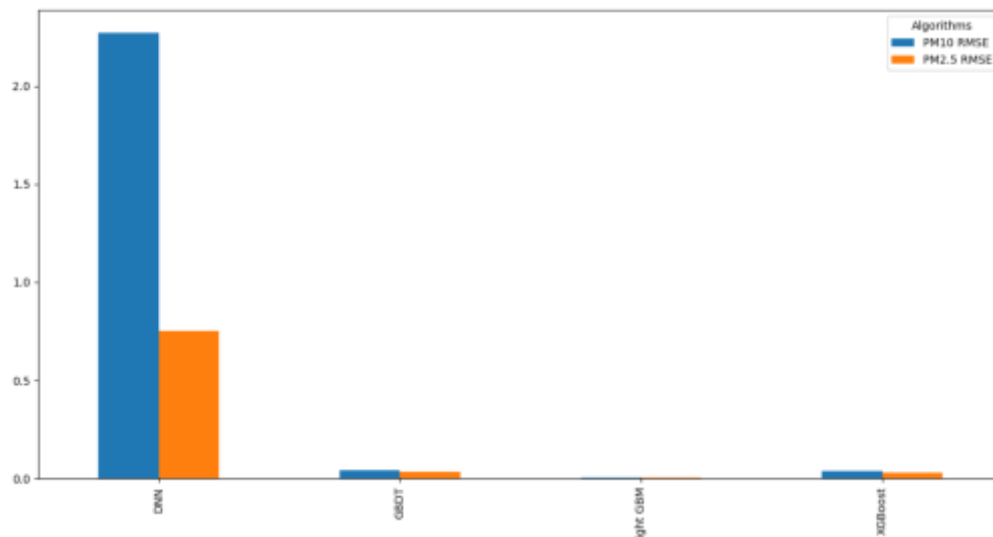Figure 8: PM10 Air Quality Prediction performance evaluation.



Figure 9: Comparison Graph of all Models.

## 5. CONCLUSION

A substantial body of research on the effects of air pollution on asthma has been published in the past 5 years, adding to the body of knowledge that has accumulated over several decades. Presently, short-term exposures to ozone, nitrogen dioxide, sulphur dioxide, PM2·5, and TRAP is thought to increase the risk of exacerbations of asthma symptoms. Increasing amounts of evidence also suggest that long-term exposures to air pollution, especially TRAP and its surrogate, nitrogen dioxide, can contribute to new-onset asthma in both children and adults. Much more about the mechanisms that are involved with exacerbations induced by pollution and onset of asthma needs to be understood, but oxidative stress and immune dysregulation are probably both involved. Young children with asthma, especially those growing up in poor neighborhoods, are at increased risk of adverse effects from exposures to air pollution. Unravelling which components of the traffic pollution mixture are responsible for asthma exacerbations and onset is a substantial challenge. Improved air quality to prevent exacerbations and new cases of asthma will require strong governmental efforts to move economies in both developed and developing countries away from combustion of fossil fuels for transportation and energy production; this approach is also needed to mitigate climate change.

## REFERENCES

[1] Lim SS, Vos T, Flaxman AD, et al. A comparative risk assessment of burden of disease and injury attributable to 67 risk factors and risk factor clusters in 21 regions, 1990–2010: a systematic analysis for the Global Burden of Disease Study 2010. Lancet. 2012; 380:2224–60.

[2] Brauer M, Amann M, Burnett RT, et al. Exposure assessment for estimation of the global burden of disease attributable to outdoor air pollution. Environ Sci Technol. 2012; 46:652–60.

[3] 3. Perez L,Declercq C, Iniguez C, et al. Chronic burden of near-roadway traffic pollution in 10 European cities (APHEKOM network). Eur Respir J. 2013; 42:594–605. [PubMed: 23520318]

[4] Wong GWK, Chow CM. Childhood asthma epidemiology: insights from comparative studies of rural and urban populations. PediatrPulmonol. 2008; 43:107–16.

[5] Robinson CL, Baumann LM, Romero K, et al. Effect of urbanisation on asthma, allergy and airways inflammation in a developing country setting. Thorax. 2011; 66:1051–57.

[6] Brunekreef B, Stewart AW, Anderson HR, Lai CKW, Strachan DP, Pearce N. Self-reported truck traffic on the street of residence and symptoms of asthma and allergic disease: a global relationship in ISAAC phase 3. Environ Health Perspect. 2009; 117:1791–98.

[7] Friedman MS, Powell KE, Hutwagner L, Graham LM, Teague WG. Impact of changes in transportation and commuting behaviors during the 1996 Summer Olympic Games in Atlanta on air quality and childhood asthma. JAMA. 2001; 285:897–905.

[8] O'Connor GT, Neas L, Vaughn B, et al. Acute respiratory health effects of air pollution on children with asthma in US inner cities. J Allergy Clin Immunol. 2008; 121:1133–39. Weinmayr G, Romeo E, De Sario M, Weiland SK, Forastiere F. Short-term effects of PM10 and NO2 on respiratory health among children with asthma or asthma-like symptoms: a systematic review and meta-analysis. Environ Health Perspect. 2010; 118:449–57.

[9] Health Effects Institute. Traffic-related air pollution: a critical review of the literature on emissions, exposure, and health effects. Boston, MA: Health Effects Institute; 2010. Panel on the Health Effects of Traffic-Related Air Pollution.

[10] Dong G-H, Chen T, Liu M-M, et al. Gender differences and effect of air pollution on asthma in children with and without allergic predisposition: northeast Chinese children health study. PLoS One. 2011; 6:e22470.

[11] Nishimura KK, Galanter JM, Roth LA, et al. Early-life air pollution and asthma risk in minority children. The GALA II and SAGE II studies. Am J Respir Crit Care Med. 2013; 188:309–18.

[12] Jacquemin B, Sunyer J, Forsberg B, et al. Home outdoor NO2 and new-onset of self-reported asthma in adults. Epidemiology. 2009; 20:119–26.