# ANALYSIS AND PREDICTION OF AIR POLLUTION USING MEMORY BASED LEARNING APPROACHES

**A PROJECT REPORT**

*Submitted by*

**JAYAKANT P [211420104108]**

**RADHEY SHYAM R [211420104210]**

**RAGUL A [211420104212]**

*in partial fulfilment for the award of the degree*
*of*

**BACHELOR OF ENGINEERING**

**in**

**COMPUTER SCIENCE AND ENGINEERING**

**PANIMALAR ENGINEERING COLLEGE**

**(An Autonomous Institution, Affiliated to Anna University, Chennai)**

**MARCH 2024**

# PANIMALAR ENGINEERING COLLEGE

### (An Autonomous Institution, Affiliated to Anna University, Chennai)

## BONAFIDE  CERTIFICATE

Certified that the project report "**ANALYSIS AND PREDICTION OF AIR POLLUTION USING  MEMORY BASED LEARNING APPROACHES**" is the bonafide work of **JAYAKANT P [211420104108], RADHEY SHYAM R [211420104210] and RAGUL A [211420104212] and** who carried out the project work under my supervision.

**Signature of the HOD with date**
**Dr L. JABASHEELA M.E., Ph.D.,**
**Professor and Head,**

Department of Computer Science and Engineering, Panimalar Engineering College Chennai - 123

**Signature of the Supervisor with date**
**Mr. D. ELANGOVAN, M.E.**
 **Associate  Professor**

Department of Computer Science and Engineering, Panimalar Engineering College Chennai - 123

Submitted for the Project Viva – Voce examination held on _____

**INTERNAL  EXAMINER**                          **EXTERNAL  EXAMINER**

# DECLARATION BY THE STUDENT

We, **JAYAKANT P [211420104108], RADHEY SHYAM R [211420104210] and RAGUL A [211420104212],** hereby declare that the project report titled "ANALYSIS AND PREDICTON OF AIR POLLUTION  USING MEMORY BASED LEARNING APPROACHES, under the guidance of Mr. D ELANGOVAN is the original work done by us and we have not plagiarized or submitted to any other degree in any university by us.

JAYAKANT P [211420104108]
RADHEY SHYAM R [211420104210]
RAGUL A [211420104212]

# ACKNOWLEDGEMENT

# ABSTRACT

Air pollution poses a significant threat to public health and environmental sustainability worldwide. Traditional methods of monitoring and mitigating air pollution often face challenges such as limited spatial coverage, high costs and delays in data processing. To address these issues, this paper explores the application of machine learning (ML) techniques to analyze and predict air pollution levels with enhanced accuracy and efficiency. It involves examining the complexities of data collection and preprocessing, emphasizing the integration of diverse data sources such as ground-based monitoring stations, satellite observations. By employing advanced feature selection and engineering methods, the paper identifies key predictors of air pollution including traffic patterns, industrial emissions and many more.

Various ML models including regression algorithms, decision trees and neural networks are systematically evaluated to determine their suitability for different prediction tasks. The evaluation process considers factors such as model accuracy, interpretability and computational efficiency. The models capable of capturing the intricate spatial and temporal dynamics of air pollution usingrigorous validation techniques, the effectiveness of the proposed ML models is demonstrated, showcasing their potential to enhance air quality monitoring and forecasting. The study also discusses the integration of these findings into decision support systems for air quality management, highlighting the benefits for policymakers, environmental agencies and the general public.

In summary, this paper contributes to the growing field of environmental informatics by offering insights into the development of more effective, data-driven strategies for air pollution mitigation and public health protection. By leveraging the power of ML techniques, this research aims to address the challenges associated with traditional air quality monitoring methods and pave the way for a more sustainable and healthier future.

# TABLE OF CONTENTS

# LIST OF SYMBOLS

| SYMBOL | Name | FUNCTION |
|---|---|---|
| Actor figure | Actor | Represents an external entity, such as a user or system, interacting with the system being modeled. |
| Object Lifeline | Object Lifeline | Represents the existence timeline of an object in a sequence diagram. |
| 1: [condition] message name → | Message | communication or interaction between objects in a sequence diagram. |
| Class Object | Object | Represents an instance of a class in a diagram, encapsulating data and behavior. |
| Use Case | Use Case | Describes a system's functionality or a specific action that can be performed, often initiated by an actor, in a use case diagram. |
| <<include>> <<extend>> | Relationships | Depict associations, dependencies, or connections between UML elements, such as classes, in diagrams. |

# CHAPTER 1

# INTRODUCTION

Air pollution is a critical environmental issue affecting numerous countries aroundthe world and India is no exception. The alarming levels of air pollution in many Indian cities have serious implications for public health, the environment and the overall quality of life. This work presents an analysis of the current state of air pollution in India and a prediction model to estimate future pollution levels. The analysis of air pollution in India involves examining various factors such as industrial emissions, vehicular pollution, biomass burning and dust particles. Additionally, meteorological conditions, including temperature, wind speed and rainfall also play a significant role in air pollution levels. The data collected from monitoring stations, satellite imagery and other relevant sources are used for analysis. The study utilizes advanced data analytics techniques including machine learning algorithms to develop a predictive model for estimating air pollution levels. The historical pollution data along with meteorological parameters are used as input to train the model.

## 1.1 PROBLEM DEFINITION

Air pollution is a critical environmental issue affecting numerous countries around the world and India is no exception. The alarming levels of air pollution in many Indian cities have serious implications for public health, the environment and the overall quality of life. This abstract presents an analysis of the current state of air pollution in India and a prediction model to estimate future pollution levels. The analysis of air pollution in India involves examining various factors such as industrialemissions, vehicular pollution, biomass burning and dust particles. Additionally, meteorological conditions including temperature, wind speed and rainfall also play

a significant role in air pollution levels. Data from monitoring stations, satellite imagery and other relevant sources are used to gather information for the analysis. The study utilizes advanced data analytics techniques, including machine learning algorithms, to develop a predictive model for estimating air pollution levels. Historical pollution data, along with meteorological parameters, are used as inputs totrain the model. The model's objective is to forecast pollution levels in different regions of India for specific timeframes, such as daily, weekly, or monthly intervals. The abstract concludes by discussing the potential applications of the analysis and prediction model. It highlights the significance of such models in aiding pAir pollution remains a critical global challenge, directly impacting public health, environmental sustainability, and economic structures. The intricate nature of air pollution necessitates advanced methodologies for monitoring, analysis, and prediction to effectively manage and mitigate its adverse effects. Traditional air quality monitoring systems, while essential, face limitations regarding spatial coverage, cost, and real-time data processing capabilities. Consequently, there's a pressing need for innovative approaches that can complement existing systems and provide deeper insights into air quality dynamics. This paper proposes leveraging machine learning (ML) techniques as a potent solution to these challenges, focusing on the prediction and analysis of air pollution levels with enhanced accuracy and efficiency.

The crux of employing ML in air pollution studies lies in the meticulous collection, preprocessing, and analysis of data, which forms the foundation for any predictive modeling endeavor. Air quality datasets are inherently complex, sourced from a variety of monitoring networks, including ground stations, satellites, and increasingly, Internet of Things (IoT) sensors. Effective feature selection is pivotal, as it involves identifying and extracting the most relevant variables influencing air pollution, such as meteorological conditions, traffic patterns, industrial emissions, and more. This process is not only critical for model performance but also for

understanding the underlying factors contributing to air pollution.

Choosing the right ML model is another critical aspect of this research. The model must not only capture the complex relationships between various predictors and air pollution levels but also be interpretable and scalable. Various ML algorithms, including regression models, tree-based methods, and neural networks, will be evaluated to identify the most suitable approach for different types of air quality prediction tasks, such as short-term forecasting and long-term trend analysis. The evaluation of these models through rigorous validation techniques is essential to ensure their reliability and effectiveness in real-world scenarios.

Moreover, air pollution is inherently influenced by spatial and temporal dynamics, necessitating models that can accurately capture these variations. Advanced ML models that can handle spatial-temporal data are crucial for making accurate predictions across different locations and times. Integrating these models into decision support systems can significantly benefit policymakers, environmental agencies, and the public, providing actionable insights for air quality management and pollution control strategies.

In summary, this paper aims to address the multifaceted challenge of analyzing and predicting air pollution through advanced machine learning techniques. By tackling the issues of data collection and preprocessing, feature selection and engineering, model selection and evaluation, and capturing spatial-temporal dynamics, this research seeks to provide a comprehensive framework for enhancing air quality monitoring and prediction efforts. The ultimate goal is to contribute to the development of more effective, data-driven solutions for mitigating air pollution and protecting public health and the environment.

## 1.2 EXISTING SYSTEM

Accurate and comprehensive air pollution data is crucial for effective environmental management and public health protection. However, missing data in air pollution monitoring datasets can hinder the ability to analyse and understand pollution patterns. This abstract presents a novel approach for hierarchical recovery of missing air pollution data using an improved Long-Short Term Context Encoder (LSTCe) network. The proposed method leverages thehierarchical structure of air pollution data, where various pollutants are measured at different monitoring stations across a region. The LSTCe network is designedto capture long-term and short-term contextual information from neighbouring stations and time intervals, respectively. By exploiting the spatial and temporal dependencies, the network can effectively recover missing data points. To enhancethe performance of the LSTCe network, several improvements are introduced. These include the incorporation of attention mechanisms to focus on relevant features and the utilization of residual connections to alleviate information loss during network training. Additionally, data augmentation techniques are appliedto address data sparsity and improve the network's generalization ability

The existing predictive models have several limitations, including the need for extensive historical data for accurate predictions, difficulty in capturing nonlinear relationships between different pollutants and contributing factors and the inability to generalize well across different regions without retraining the models with localized data. Moreover, real-time prediction capabilities are often limited, reducing the effectiveness of these models in guiding immediate policy and public health responses.

**Disadvantages**

- Complexity in training the data model.

- Low accuracy predictions.

- undeployed model leading to instability.

- Relies only on the implemented algorithm for output.

## 1.3 PROPOSED SYSTEM

Air pollution poses a critical environmental challenge in India, significantly impacting public health and the overall quality of life. The effective mitigation strategies and policy-making rely on accurate analysis and prediction of air pollution levels. To address this need, a proposed system is proposed to leverage memory-based learning techniques for the analysis and prediction of air pollution in India.

The system utilizes a comprehensive set of inputs including pollutant concentrations, meteorological variables, geographical features and temporal information to train machine learning models. By incorporating such diverse data sources, the system can capture complex relationships and patterns inherent in air pollution dynamics. This holistic approach enables the system to provide more accurate and nuanced predictions compared to traditional methods.

The proposed system offers several potential applications in air quality management. It can aid in designing targeted pollution control measures by identifying key factors driving pollution levels and assessing their impact. Then the system issues timely alerts based on predictive models, enabling the authorities to take proactive measures for mitigating pollution spikes and protect public health as well as it is capable of generating all the possible diseases caused

due to the predicted toxic levels from the gas sample which can be used by the government organizations to analyze and prepare for preventing and spreading awareness to people thereby protecting public health and environmental balance. By optimizing resource allocation based on forecasted pollution levels, the system can enhance the efficiency of air quality management efforts in India.

In summary, this system presents a novel approach to tackle the challenges of air pollution in India by applying machine learning algorithms. By leveraging memory-based learning techniques and integrating diverse data sources, the proposed system holds promise in mitigating the adverse effects of air pollution and improving public health outcomes.

**Advantages:**

- Machine learning techniques such as memory-based approaches are used to build a statistical predictive model.

- More than two algorithms are implemented for comparative analysis.

- A Full-stack application for the model is deployed which can be easily accessed.

- Enhanced accuracy & performance level.

## 1.4 INTRODUCTION TO MEMORY-BASED MAHCINE LEARNING TECHNIQUES

Memory-based learning approaches, also known as instance-based learning or lazy learning, constitute a subset of machine learning techniques that rely on stored instances of training data to make predictions or classifications. Unlike traditional models that require explicit generalization through parameter estimation, memory-based methods defer learning until the prediction phase,

where they retrieve and analyze relevant instances from the training dataset.

The core principle behind memory-based learning is the notion that similar instances in the training data will exhibit similar behaviors or outcomes. Therefore, instead of abstracting patterns from the entire dataset during training, memory-based models directly compare the input instance with the stored training instances to make predictions. This approach offers several advantages such as simplicity, flexibility and the ability to adapt changing data distributions.

One of the key components of memory-based learning is the similarity measure used to compare instances. The common employed similarity metrics include Euclidean distance, cosine similarity and Pearson correlation coefficient among others. By quantifying the similarity between instances, memory-based models can effectively identify relevant training examples for making predictions. The memory-based learning techniques encompass various algorithms with k- nearest neighbors (KNN) being one of the most prominent examples. In KNN, predictions are made based on the majority vote or weighted average of the labels of the k nearest neighbors to the input instance in the feature space. Other memory-based methods are locally weighted regression, case-based reasoning and collaborative filtering. While memory-based learning approaches offer simplicity and interpretability, they also exhibit limitations. Their reliance on stored training instances can lead to computational inefficiency, especially for large datasets. Additionally, memory-based models may suffer from the curse of dimensionality and struggle to generalize well to unseen data if the training instances are not representative of the underlying data distribution.

Despite these challenges, memory-based learning approaches remain valuable tools in machine learning, particularly in scenarios where interpretability and adaptability are prioritized over computational efficiency. Their ability to

leverage stored instances of training data to make predictions makes them well-suited for tasks such as classification, regression and recommendation systems contributing to the diverse landscape of machine learning techniques.

# CHAPTER 2
# LITERATURE SURVEY

**[1] Sai Bhargav Kasetty, S. Nagini. A Survey Paper on an IoT-based Machine Learning Model to Predict Air Pollution Levels.**

In today's quickly developing world, air pollution is a major worry. As pollution rises in the earth's atmosphere every day, so do its rates. The problem is typically caused by toxic fumes released into the air as a result of a broad range of human activities. It is necessary for the current generation to understand the risks posed by these circumstances and to take effective measures to minimize them before the ecosystem is destroyed. This study examines studies of machine learning approaches (MLAs or MLTs) and Internet of Things in predicting air quality, as well as forecasting air pollution levels so individuals can take action to reduce pollution.

This paper focuses on the escalating issue of air pollution, primarily driven by a wide range of human activities leading to the release of toxic fumes. On recognizing the urgency to address these risks, the paper examines the use of Machine Learning Approaches (MLAs or MLTs) and the Internet of Things (IoT) in predicting and forecasting air quality. This enables individuals and authorities to take proactive measures to mitigate pollution levels. The integration of IoT with machine learning offers a robust framework for real-time data collection and analysis,which is vital for accurate air pollution forecasting.

**Advantages**: The combination of IoT and MLAs for predicting air pollution represents a cutting-edge approach, leveraging real-time data for accurate forecasting. By forecasting air pollution levels, individuals and authorities can take timely actions to reduce pollution, potentially leading to improved public health and environmental protection. The IoT devices provide real-time data, enhancing the accuracy of ML-based predictions. This real time aspect is critical in dealing with dynamic environmental factors.

**Disadvantages**: The use of IoT devices raises concerns about data privacy and

security, as these devices often collect large amount of sensitive data. By setting up a network of IoT devices and maintaining the ML model can be costly, potentially limiting its application in resource-constrained settings. By integrating diverse IoT devices and ensuring the seamless functioning of the ML model can be complex and require technical expertise.

**[2] Peijiang Zhao, Koji Zettsu "Convolution Recurrent Neural Networks Based Dynamic Transboundary Air Pollution Predictiona 2019"**

Transboundary air pollution is one of the main sources of air pollution in island cities. However, the transboundary pollution confounded by local emission, meteorological conditions and it is difficult to predict. At present, most of urban air pollution prediction methods do not predict with transboundary air pollution. Therefore, we introduce a dynamic transboundary air pollution prediction approach based on convolutional recurrent neural networks(D-CRNN) which: (i) Divides the prediction inputs into prediction locations and transboundary air pollution sources (ii) Using two different convolutional recurrent neural networks to solve the spatial-temporal feature of each inputs. (iii) Through a transboundary prediction network to integrate the spatial-temporal feature of prediction locations with the spatial-temporal feature of transboundary air pollution sources in a dynamic asynchronous method. Then use those mixed features to predict the air pollution. To evaluate DCR.NN model with the local atmospheric monitoring data in Kyushu, Japan and the transboundary air pollution data from 33 coastal cities in eastern Asia from January 2015 to July 2017. The results show that D-CRNN model has achieved 86.2%, 78.6% accuracy of total prediction and transboundary air pollution in next 6 hours.

**Advantages**: CRNNs combine the strengths of Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), making them highly effective in processing and predicting spatial-temporal data, which is crucial in air pollution prediction. The model's ability to predict transboundary air pollution is significant,

considering that air pollution is not confined to geographic boundaries. This can aid in international policy-making and collaborative efforts for pollution control. CRNNs are adept at handling real-time data, providing timely predictions that can be critical for taking immediate measures to mitigate pollution impacts.

**Disadvantages**: CRNNs can be computationally intensive due to their complex architecture, requiring significant processing power and potentially leading to higher operational costs. The accuracy of the predictions is highly dependent on the quality and quantity of the input data. The poor quality data can lead to inaccurate predictions. By implementing and maintaining CRNN models, it requires a high level of technical expertise in machine learning and environmental science which might be a barrier for some organizations.

**[3] Shweta Taneja, Nidhi Sharma, Kettun Oberoi, Yash Navoria "Predicting trends in air pollution in Delhi using data mining 2018"**

Over the past years the development and urbanization in Delhi has led to increase in air pollution. This has led to study and research in this area. We have used data mining to analyze the existing trends in air pollution in Delhi and make prediction about the future. The data mining techniques used are linear regression andmultilayer perceptron. We have seen the trends of various air pollutants like sulphur dioxide($SO_2$), nitrogen dioxide ($NO_2$), particulate matter (PM), carbon monoxide (CO), ozone ($O_3$). By using the above techniques we have observed that there will be an increase in amount of PM 10 by 45.9% in coming years. However amount of CO and NO2 may show slight increase due to increasing number of 2 wheelers on road. The other pollutants like $SO_2$ may show decrease due to usage of non sulphur fuel and stringent pollution control measures.

**Advantages**: Predicting trends in air pollution using data mining techniques can serve as an early warning system, enabling authorities to take proactive measures to

mitigate pollution levels before they reach hazardous levels. Data mining can help identify patterns and factors contributing to air pollution trends, allowing for more efficient allocation of resources and targeted interventions to reduce pollution sources. Insights derived from data mining can inform policy formulation and regulatory measures aimed at controlling air pollution, leading to more effective and evidence-based environmental policies.

**Disadvantages**: Data mining for predicting air pollution trends relies on the availability of high-quality and comprehensive datasets. However, data on air pollution may be limited in coverage or subject to measurement errors, affecting the accuracy of predictions. Data mining models used for predicting air pollution trends can be complex and difficult to interpret, especially for stakeholders without a background in data science. This may hinder the adoption and implementation of predictive models in decision-making processes. Predictive models developed using data mining techniques may have inherent uncertainties and limitations, particularly when dealing with complex environmental systems. The uncertainties in model predictions can undermine confidence in the reliability of the results.

**[4] Shweta Taneja, Nidhi Sharma ,Kettun Oberoi ,Yash Navoria "Predicting trends in air pollution in Delhi using data mining 2019"**

Over the past years the development and urbanization in Delhi has led to increase in air pollution. This has led to study and research in this area. We have used data mining to analyze the existing trends in air pollution in Delhi and make prediction about the future. The data mining techniques used are linear regression and multilayer perceptron. We have seen the trends of various air pollutants like sulphur dioxide($SO_2$), nitrogen dioxide ($NO_2$), particulate matter (PM), carbon monoxide (CO), ozone ($O_3$). By using the above technique we have observed that there will be an increase in amount of PM 10 by 45.9% in coming years, however amount of CO and NO2 may show slight increase due to increasing number of 2 wheelers on road. The other

pollutants like SO 2 may show decrease due to usage of non sulphur fuel and stringent pollution control measures.

**Advantages**: Data mining techniques are capable of efficiently processing and extracting valuable insights from large datasets, which is essential for analyzing complex air pollution data. These techniques excel at identifying patterns and trends in historical pollution data, which can be crucial for predicting future air pollution levels. By leveraging historical data, data mining can potentially provide more accurate predictions of pollution trends compared to traditional statistical methods. Data mining allows for the development of customized models that can be tailored to the specific environmental and pollution variables relevant to Delhi.

**Disadvantages**: The accuracy of predictions is heavily reliant on the quality and completeness of the input data. Poor data quality can lead to unreliable predictions. Some data mining models, particularly advanced ones like neural networks, can be complex and difficult to interpret, making it challenging to understand how predictions are made. Processing large datasets with sophisticated data mining techniques can be computationally intensive, requiring significant computational resources. Data mining models, especially those that are not properly regularized, can overfit to the historical data, leading to poor performance in predicting future trends.

**[5] Shahan Salim, Irfhana Zakir Hussain, Jasleen Kaur, Plinio P. Morita "An Early Warning System for Air Pollution Surveillance: An IoT Based Big Data Framework to Monitor Risks Associated with Air Pollution 2020"**
Air pollution is a major public health issue that can have far-reaching consequences for human health. Accurately quantifying its impact and effects can be a challenging task. However, with the advent of the Internet of Things (IoT) and big data technologies, real-time monitoring of air pollution levels is now possible, enabling prompt and data-driven measures to mitigate its negative effects. This paper proposes

the development of an agnostic ecosystem that collects big data from various sensors and analyzes and predicts harm using state-of-the-art AI and deep learning techniques. The ecosystem can provide public health officials and researchers with the necessary tools to monitor air pollution levels and design effective policies to lessen the harmful effects of air pollution on human health.

**Advantages**: The data-driven approach facilitates the identification of trends and patterns in air pollution levels over time. By leveraging data mining algorithms, such as regression analysis or time series forecasting, it becomes possible to develop predictive models that can forecast future air pollution trends. These models can assist policymakers and environmental agencies in implementing proactive measures to mitigate pollution risks. The implementation of an Early Warning System (EWS) based on data mining techniques enables real-time monitoring of air pollution levels. The integration of Internet of Things (IoT) devices in the framework enhances data collection capabilities, enabling continuous monitoring of air quality parameters such as particulate matter, ozone levels, and nitrogen dioxide concentrations. This real-time data collection ensures the timely detection of pollution events.

**Disadvantages**: The accuracy and reliability of air pollution data collected by IoT devices may vary due to factors such as sensor calibration, environmental conditions, and data transmission errors. Ensuring the quality and reliability of the data is crucial for the effectiveness of the predictive models and early warning system. The processing of large volumes of data generated by IoT devices and other sources requires substantial computational resources and advanced data processing techniques. Managing and analyzing big data in real-time can pose challenges in terms of scalability and performance. Implementing predictive modeling algorithms for air pollution trend prediction involves dealing with complex mathematical modelsand computational algorithms. Understanding and optimizing these algorithms require specialized knowledge and expertise in data mining and machine learning.

**[6] G. Jignesh Chowdary, Suganya. G, Premalatha. M proposed a paper "Effective prediction of cardiovascular disease using cluster of machine learning algorithms"** which depicts that Cardiovascular diseases are one of the diseases that account for the loss of millions of lives each year. Lack of early prediction is the primary reason for the loss of lives, and this encourages researchers to develop intelligent systems for better prediction. In this paper, a novel ensemble methodology is introduced which uses the voting of Logistic Regression (LR), Random Forest (RF), Artificial Neural Network activated with ReLU function (NNR), K-Nearest Neighbours (KNN) and Gaussian Naive Bayes (GNB) to predict the possibility of heart disease. The model is developed using Python-based Jupyter Notebook and Flask and is trained using the standard dataset from Kaggle. The model is tested and evaluated based on accuracy, precision, specificity, sensitivity, error. Testing witnessed an accuracy of 89% and a precision of 91.6%, along with a sensitivity of 86% and specificity of 91%. The results upon comparison with the individual models witness the better accuracy of using ensemble modelling and hence a better prediction leading to life-saving.

**Advantages**: The ensemble methodology proposed in the paper achieves an accuracy of 89%, which is crucial for effectively predicting cardiovascular diseases (CVD) and enabling timely interventions. With a precision of 91.6%, the model ensures a high percentage of correct positive predictions among all positive predictions, reducing false positives and unnecessary interventions.

**Disadvantages**: The performance of the model is contingent upon the quality and representativeness of the dataset from Kaggle, which may limit its applicability to real-world clinical settings with diverse patient populations. The ensemble model, comprising multiple algorithms, may pose challenges in interpreting the underlying decision-making process, potentially hindering clinical adoption and trust among

healthcare professionals.

**[7] R. Chitra1 and V. Seenivasagam,** "**Review of Heart Disease Prediction System Using Data Mining and Hybrid Intelligent Techniques 2020**"**,** stated that the Healthcare industry generally clinical diagnosis is done mostly by doctor's expertise and experience. Computer Aided Decision Support System plays a major role in medical field. With the growing research on heart disease predicting system, it has become important to categories the research outcomes and provides readers with an overview of the existing heart disease prediction techniques in each category.Neural Networks are one of many data mining analytical tools that can be utilized to make predictions for medical data. From the study it is observed that Hybrid Intelligent Algorithm improves the accuracy of the heart disease prediction system. The commonly used techniques for Heart Disease Prediction and their complexities are summarized in this paper

**Advantages**: Integration of data mining and hybrid intelligent techniques enhances the accuracy of heart disease prediction compared to traditional diagnostic methods, potentially leading to more precise and timely interventions. The paper provides a comprehensive overview of existing heart disease prediction techniques, categorizing research outcomes and summarizing commonly used methods and their complexities. This facilitates a deeper understanding of the current state-of-the-art in heart disease prediction.

**Disadvantages**: Implementing data mining and hybrid intelligent techniques can be complex and resource-intensive, requiring specialized knowledge and significant computational resources, which may limit their widespread adoption, especially in resource-constrained healthcare settings. Risk of Over-Reliance and Lack of Personalization: There's a risk of over-reliance on algorithmic predictions, potentially leading to a lack of personalized treatment tailored to individual patient needs. Human

judgment and expertise might be undervalued, resulting in less nuanced decision-making.

**[8] Avinash Golande, Pavan Kumar T, Heart Disease Prediction Using Effective Machine Learning Techniques** paper states that in today's era deaths due to heart disease has become a major issue approximately one person dies per minute due to heart disease. This is considering both male and female category and this ratio may vary according to the region also this ratio is considered for the people of age group 25-69. This does not indicate that the people with other age group will not be affected by heart diseases. This problem may start in early age group also and predict the cause and disease is a major challenge nowadays. Here in this paper, we have discussed various algorithms and tools used for prediction of heart diseases.

**Advantages**: By utilizing effective machine learning techniques, such as various algorithms discussed in the paper, the heart disease prediction system can achieve higher accuracy compared to traditional methods. This can lead to early detection and intervention, potentially saving lives. Early Detection Across Age Groups: The paper acknowledges that heart diseases can affect individuals across various age groups, not just limited to the 25-69 age bracket. By employing machine learning algorithms, the system can potentially identify risk factors and symptoms early, even in younger individuals, allowing for proactive management of heart health.

**Disadvantages**: The effectiveness of machine learning models for heart disease prediction relies heavily on the quality and representativeness of the underlying data. Inaccurate or biased data can lead to erroneous predictions and compromised system performance. Interpretability and Transparency Challenges: Some machine learning algorithms, particularly complex models like neural networks, may lack interpretability, making it challenging for healthcare professionals to understand the underlying factors driving the predictions. This lack of transparency could hinder trust

17

and acceptance of the prediction system among clinicians and patients.

**[9] Animesh Hazra, 2Subrata Kumar Mandal, Heart Disease Diagnosis and Prediction Using Machine Learning and Data Mining Techniques**: A Review

A popular saying goes that we are living in an "information age". Terabytes of data are produced every day. Data mining is the process which turns a collection of data into knowledge. The health care industry generates a huge amount of data daily. However, most of it is not effectively used. Efficient tools to extract knowledge from these databases for clinical detection of diseases or other purposes are not much prevalent. The aim of this paper is to summarize some of the current research on predicting heart diseases using data mining techniques, analyse the various combinations of mining algorithms used and conclude which technique(s) are effective and efficient. Also, some future directions on prediction systems have been addressed.

**Advantages**: The paper highlights the importance of data mining techniques in extracting knowledge from the vast amount of healthcare data generated daily. This knowledge can be instrumental in improving heart disease diagnosis and prediction, leading to better patient outcomes. Efficient Clinical Detection: By summarizing current research on predicting heart diseases using data mining techniques, the paper aims to contribute to more efficient clinical detection of heart diseases. This can potentially lead to early intervention and improved management of cardiovascular health.

**Disadvantages**: The effectiveness of data mining techniques relies heavily on the quality and availability of healthcare data. Issues such as incomplete or inaccurate data, as well as data fragmentation across different healthcare systems, can hinder the accuracy and reliability of predictions. Analyzing various combinations of mining algorithms may introduce complexity, making it challenging to determine the optimal approach for heart disease prediction. This complexity can increase the computational

burden and require expertise in both data mining and healthcare domain knowledge.

**[10] Mangesh Limbitote, Dnyaneshwari Mahajan A Survey on Prediction Techniques of Heart Disease using Machine Learning 2020.**

Heart is one of the most important part of the body. It helps to purify and circulate blood to all parts of the body. Most number of deaths in the world are due to Heart Diseases. Some symptoms like chest pain, faster heartbeat, discomfort in breathing are recorded. This data is analyzed on regular basis. In this review, an overview of the heart disease and its current procedures is firstly introduced. Furthermore, an in- depth analysis of the most relevant machine learning techniques available on the literature for heart disease prediction is briefly elaborated. The discussed machine learning algorithms are Decision Tree, SVM, ANN, Naive Bayes, Random Forest, KNN. The algorithms are compared on the basis of features. We are working on the algorithm with best accuracy. This will help the doctors to assist the heart problem easily.

**Advantages**: The paper provides an overview of heart disease, highlighting its significance as one of the leading causes of death worldwide. This comprehensive understanding sets the context for exploring prediction techniques using machine learning. Insight into Current Procedures: By discussing current procedures for analyzing heart disease data, the paper offers insights into existing methodologies used in healthcare settings. Understanding these procedures is essential for evaluatingthe effectiveness and efficiency of machine learning techniques in comparison.

**Disadvantages**: While the paper discusses several well-established machine learning algorithms for heart disease prediction, it may lack exploration of newer or emerging techniques that could potentially offer improved accuracy and efficiency. Generalization of Results: The comparison of machine learning algorithms may not account for the variability in datasets and clinical settings, potentially leading to overgeneralization of results. Different datasets and patient populations may exhibit unique characteristics that impact the performance of algorithms differently.

# CHAPTER 3

# THEORETICAL BACKGROUND

## 3.1 IMPLEMENTATION ENVIRONMENT

### 3.1.1  Hardware Requirements:

- Processor: Intel i3

- Hard disk: minimum 80 GB

- RAM: minimum 2 GB

### 3.1.2  Software Requirements:

- Operating system: Windows 10 or higher

- Tool: Any IDE

For implementing a project aimed at analyzing and predicting air pollution in India using memory-based learning approaches, you would need to consider several key components in your environment setup. This setup spans data collection, processing, model development, and deployment phases. Below is an outline of what such an environment might look like, covering software, hardware, and datasets.

### 3.1.3  Software and Tools

#### 3.1.3.1  Programming Languages:

Python: Widely used in data science and machine learning projects for its simplicity

and the vast array of libraries available.

### 3.1.3.2  Machine Learning Libraries:

Scikit-learn: Provides simple and efficient tools for data mining and data analysis. It's built on NumPy, SciPy, and matplotlib.

TensorFlow or PyTorch: For more complex models, including deep learning networks that might be useful for analyzing complex patterns in air quality data.

Pandas: Essential for data manipulation and analysis.

NumPy: Adds support for large, multi-dimensional arrays and matrices, along with a large collection of high-level mathematical functions to operate on these arrays.

### 3.1.3.3  Data Visualization Tools:

Matplotlib and Seaborn: For creating static, animated, and interactive visualizations in Python.

Tableau or Power BI (optional): For more advanced interactive visualizations, especially useful in presentations or reports for non-technical stakeholders.

### 3.1.3.4  Development Environment:

Jupyter Notebooks or Google Colab: Ideal for interactive development and documentation of data analysis and machine learning models.

### 3.1.3.5  Hardware

Depending on the scale of data and complexity of the models, you might need:

A decent CPU (Intel i3/i5/Ryzen 5 or 7) for basic data processing and machine

learning tasks.

A high-end GPU (e.g., NVIDIA GTX 1080 or better) for deep learning models. Cloud services like AWS, Google Cloud, or Azure offer scalable GPU resources.

Sufficient RAM (2GB or more) for handling large datasets in-memory.

### 3.1.3.6 Datasets

Air Quality Data:

Government and non-governmental organizations' databases, such as the Central Pollution Control Board (CPCB) or World Air Quality Index project, for historical air quality data.

Satellite data for broader geographical coverage.

### 3.1.3.7 Auxiliary Data:

Meteorological data: Temperature, humidity, wind speed, and direction can influence air pollutant levels.

Emissions data from industrial, transport, and residential sources.

Geographic and demographic data for analyzing population exposure to pollution.

### 3.1.3.8 Deployment Environment

For deploying models into production:

Flask or Django for creating a web application to display predictions.

Docker for containerization and easy deployment.

Heroku, AWS, or Google Cloud Platform for hosting the application.

This setup provides a comprehensive environment for developing a system to analyze and predict air pollution levels using memory-based learning approaches. Tailor the environment based on the specific needs, scale of your project, and available resources.

## 3.2 System Architecture



### 3.2.1 Website:

The website serves as the user interface, allowing stakeholders to interact withthe system. Users can input queries, access visualizations, and view predictions.It provides a convenient platform for users to access and utilize the system's capabilities without needing to directly interact with the underlying data and algorithms.

### 3.2.2 Dataset:

The dataset serves as the foundation of the system, containing historical data on air pollution levels in India.

It includes information on pollutant concentrations, meteorological variables (such as temperature, humidity, wind speed), geographical features (such aslocation, altitude), and temporal information (such as time of measurement).

The dataset is essential for training machine learning models and making predictions about future air pollution levels.

### 3.2.3 Preprocessing:

The preprocessing stage involves cleaning, transforming, and preparing the dataset for analysis. Tasks performed during preprocessing may includehandling missing or erroneous data, normalizing features to a consistent scale,and encoding categorical variables into numerical representations.

Preprocessing ensures that the data is suitable for analysis and helps improve the performance of machine learning algorithms.

### 3.2.4 Visualization:

Visualization techniques are employed to provide insights into the dataset andaid in understanding air pollution trends and correlations.

Visual representations such as charts, graphs, maps, and heatmaps are created to visualize pollutant concentrations over time, spatial distributions of pollutants, and relationships between pollution levels and other variables.

Visualization helps stakeholders identify patterns, anomalies, and areas of concern within the data.

### 3.2.5 Algorithm Implementation:

Memory-based learning algorithms, such as k-nearest neighbors (KNN),

collaborative filtering, or locally weighted regression, are implemented to analyze the preprocessed data and make predictions.

These algorithms utilize stored instances of training data to make predictions or classifications based on the similarity between input instances and historical data.

Memory-based learning approaches are chosen for their ability to capture complex relationships and patterns in the data without the need for explicit parameter estimation during training.

### 3.2.6 Accuracy Evaluation:

The performance of the implemented algorithms is evaluated using appropriate metrics and techniques.

Evaluation metrics may include accuracy, precision, recall, F1-score, and mean squared error, depending on the nature of the prediction task (classification or regression).

Accuracy evaluation ensures that the models provide reliable and accurate predictions and helps identify areas for improvement.

### 3.2.7 Model:

Once the algorithms are trained and evaluated, a final model is created.

This model encapsulates the knowledge learned from the data and can be usedfor making real-time predictions about air pollution levels.

The model represents the culmination of the analysis process and serves asthe core component of the predictive system.

### 3.2.8 Deployment Model:

The final model is deployed into production, allowing stakeholders to access

its predictions through the website or other interfaces.

Deployment may involve hosting the model on cloud platforms, integrating itinto existing air quality management systems, or embedding it within the website's backend infrastructure.

The deployed model enables users to obtain real-time predictions and insightsinto air pollution levels, supporting informed decision-making and proactive measures to mitigate pollution-related risks.

In summary, the system architecture encompasses a cohesive set of components designed to analyze and predict air pollution levels in India usingmemory-based learning approaches. From data collection and preprocessing to model deployment and prediction, each component plays a crucial role in providing actionable insights and supporting effective air quality managementstrategies.

.

## 3.3 Proposed Methodology

The proposed methodology for developing a system to analyze and predict air pollution in India using memory-based learning approaches, integrated into a full-stack web application, is comprehensive and multi-faceted. It includes the following stages:

### 3.3.1. Data Collection and Preprocessing:

Collection: Gather air quality data from multiple sources, including government databases, API endpoints (e.g., AQICN), and historical records.

Preprocessing: Clean and normalize the data to ensure consistency. This involves handling missing values, removing outliers, and transforming data into a format suitable for analysis.

### 3.3.2. Development of Web Application Framework:

Front-End Development: Use HTML, CSS, Bootstrap, and JavaScript to create a user-friendly interface that allows users to interact with the system, input data, and view predictions and visualizations.

Back-End Development: Utilize Django, a high-level Python web framework, to handle server-side logic, including data processing, interaction with the machine learning models, and database management.

### 3.3.3. Machine Learning Model Implementation:

Selection of Algorithms: Implement several memory-based learning algorithms such as K-Nearest Neighbors (KNN), Decision Tree Classifier, Random Forest Classifier, and Extra Trees Classifier.

Model Training and Evaluation: Train the models on preprocessed datasets, evaluating their performance using metrics like accuracy, recall, and F1-score to select the best-performing model for deployment.

### 3.3.4. Integration of Machine Learning Models with the Web Application:

Model Deployment: Integrate the trained models into the web application, allowing for real-time predictions based on user input.

Visualization and User Interaction: Develop interactive visualizations of air quality data and predictions using JavaScript libraries, enhancing the user experience.

### 3.3.5. Testing, Evaluation, and Deployment:

Testing: Perform thorough testing of both the web application and machine learning models to ensure accuracy, efficiency, and user-friendliness.

Deployment: Deploy the application on a suitable cloud platform, ensuring scalability and accessibility for users.

### 3.3.6. Continuous Improvement and Expansion:

Feedback Loop: Implement a mechanism to gather user feedback on predictions and usability, using this input to continuously improve the system.

Expansion: Plan for the future expansion of the system, including the addition of more

sophisticated machine learning models, broader datasets, and enhanced features based on user needs and feedback.

## 3.4 Input Design

### Step 1: Login/Registration

Users can input their username and password for login, while new users have the option to register by providing a username, password, and email. This ensures secure access to the system, with robust authentication methods safeguarding user accounts.

### Step 2: Home page

The home page serves as the central hub, users can view the profile, logout using the shorthand buttons, navigate to the results section where all the records are saved.

### Step 3: Providing Inputs

The user is supposed to provide the input details by selecting the appropriate variables that are analyzed from the gas mixture by selecting the molecules tab and providing the accurate values in suitable measurements.

### Step 4: Result tab

Once the user successfully provides valid input to the placeholders, he/she can click the run button to submit the values and generate the output based on the trained model, once the output is generated, the result is shown as a report format depicting the amount of toxicity along with the diseases caused due to the specific amounts of gases.

### Step 5: Logout

This input ensures users can safely end their sessions, contributing to the overall security and privacy of their interactions with the application.

**Module Design**

### 3.5.1 Sequence Diagram



### 3.5.2 Data Flow Diagram

### 3.5.3 Flow Chart Diagram

### 3.5.4 Use case Diagram

# CHAPTER 4

## SYSTEM IMPLEMENTATION

## MODULE DESCRIPTION

- Data Pre-processing
- Data Analysis of Visualization
- Implementing Extra Tree Classifier Algorithm
- Implementing Random Forest Classifier Algorithm
- Implementing XG Boost Classifier Algorithm
- Deployment Using Django

## 4.1 Data Pre processing

Validation techniques in machine learning are used to get the error rate of the Machine Learning (ML) model, which can be considered as close to the true error rate of the dataset. If the data volume is large enough to be representative of the population, you may not need the validation techniques. However, in real-world scenarios, to work with samples of data that may not be a true representative of the population of given dataset. To finding the missing value, duplicate value and description of data type whether it is float variable or integer. The sample of data used to provide an unbiased evaluation of a model fit on the training dataset while tuning model hyper parameters. In the context of our study on analysing and predicting air pollution using machine learning techniques, the data preprocessing module plays a crucial role in ensuring the reliability and accuracy of our predictive models. Given the complexity and variability of air quality data, which includes inputs from diverse sources such as satellite observations, ground-based monitoring stations, and Internet of Things (IoT) sensors, the preprocessing module is designed to address several key challenges: data cleaning, integration, normalization, and feature engineering. This section provides a detailed description of the data preprocessing module, outlining its components and

the methodologies employed to prepare the dataset for subsequent analysis and modelling.

### 4.1.1 Data Cleaning

The initial step in the preprocessing module involves data cleaning, which is essential for removing inaccuracies and inconsistencies from the data. This process includes:

**4.1.1.1 Handling Missing Values:** Missing data is a common issue in air quality datasets due to equipment malfunctions, data transmission errors, or other monitoring discrepancies. We employ various imputation techniques such as mean imputation, interpolation, and advanced ML-based imputation methods to address missing values, depending on the nature and pattern of the missing data.

**4.1.1.2 Outlier Detection and Removal**: Outliers can significantly skew the results of air quality models. We use statistical techniques and anomaly detection algorithms to identify and remove outliers, ensuring that our dataset accurately reflects typical air pollution patterns.

### 4.1.2 Data Integration

Air quality data often comes from multiple sources with varying formats, resolutions, and coverage areas. The data integration process involves:

**4.1.2.1 Harmonizing Data Formats:** Standardizing the data into a consistent format is crucial for integrating diverse data sources. This includes converting timestamps to a uniform standard, aligning measurement units, and ensuring consistent spatial resolution.

**4.1.2.2 Merging Datasets**: We merge data from different sources into a single, comprehensive dataset. This involves aligning data based on timestamps and

geographical locations, ensuring that each record provides a holistic view of air quality indicators and relevant predictors at specific times and places.

### 4.1.2.3 Data Normalization

Given the range of variables involved in air quality prediction (e.g., pollutant concentrations, meteorological conditions), normalization is critical for ensuring that all input features contribute proportionately to the predictive models. We apply normalization techniques such as Min-Max scaling and Z-score normalization to transform the data into a common scale without distorting differences in the ranges of values.

### 4.1.3 Feature Engineering

Feature engineering is pivotal in enhancing the predictive power of ML models. This involves:

**4.1.3.1 Feature Selection**: Identifying the most relevant features that impact air pollution levels. This process is informed by domain expertise, correlation analysis, and feature importance ranking methods.

**4.1.3.2 Feature Construction:** Creating new features from the existing data that might capture complex interactions or provide additional insights into air pollution dynamics. This includes calculating moving averages, lag features to capture temporal trends, and interaction terms between features.

### 4.1.4 Summary

The data preprocessing module is a foundational component of our approach to analyzing and predicting air pollution. By meticulously cleaning, integrating, normalizing, and engineering features from the collected data, we lay the groundwork for developing robust, accurate, and insightful machine learning models. This module not only enhances the quality of the input data but also ensures that our predictive

models can effectively capture the complex dynamics of air pollution, leading to more accurate forecasts and analyses.

## 4.2 Data Visualization Module Description

The Data Analysis and Visualization module in the project focuses on processing and understanding air quality data to identify trends and correlations, and presenting these insights visually. It involves cleaning and preparing data, applying statistical and machine learning techniques for deep analysis, and then using web technologies to create interactive dashboards for easy interpretation of air pollution patterns. This module helps in making informed decisions for air quality management by providing a clear view of historical and predicted pollution levels through engaging charts and maps.

### 4.2.1 Data Analysis:

The primary objectives of the data visualization module are as follows:

**4.2.1.1 Objective**: The main objective of the data analysis module is to process and analyze the air quality data to identify trends, patterns, and correlations between different pollutants and environmental factors.

**4.2.1.2 Techniques Used**: This involves using statistical methods, data mining techniques, and machine learning algorithms to understand the behavior of air pollution over time and across different geographical locations. Techniques such as time-series analysis, correlation analysis, and feature importance evaluation are employed.

**4.2.1.3 Preprocessing**: Before analysis, the data undergoes preprocessing steps such as cleaning (removing outliers and missing values), normalization (scaling data), and transformation (converting data into a suitable format for analysis).

**4.2.1.4 Feature Engineering**: This step involves creating new features from existing data to improve model performance. For air pollution data, this could include deriving

rolling averages of pollutant levels, calculating pollutant ratios, or integrating weather condition data.

### 4.2.2 Visualization:

It is a crucial component of our study on analyzing and predicting air pollution using machine learning techniques. Visualization plays a pivotal role in exploring, interpreting, and communicating insights from the air quality data, facilitating a deeper understanding of the underlying patterns and trends. In this section, we provide a detailed description of the data visualization module, outlining its objectives, methodologies, and key visualization techniques employed to enhance the analysis and interpretation of air pollution data.

**4.2.2.1 Objective**: The visualization module aims to present the analyzed data and predictions in an intuitive and interactive manner, making it easier for users to understand the air quality situation and forecasts.

**4.2.2.2 Tools and Technologies**: Utilizes web technologies like HTML, CSS, JavaScript (along with libraries such as D3.js or Chart.js) for frontend development. For backend integration, frameworks such as Django can be used to handle data processing and visualization rendering.

**4.2.2.3 Interactive Dashboards**: The module features dashboards that display air quality metrics through interactive charts, graphs, and maps. Users can select different time frames, locations, and pollutants to view specific data visualizations.

**4.2.2.4 Real-time Data Presentation**: Offers real-time air quality monitoring and prediction visualizations. This enables users to see current pollution levels and future forecasts, enhancing decision-making for public health and environmental planning.

**4.2.2.5 Insights and Patterns**: Visualizes key insights and patterns in the data, such as seasonal variations in pollutant levels, the impact of policy changes on air quality, and correlations between air pollutants and weather conditions.

**4.3 Implementing Extra tree classifier algorithm**

The module for implementing the Extra Tree Classifier algorithm in the paper focuses on utilizing this specific machine learning algorithm to predict air pollution levels in India. It involves the following key aspects:

**4.3.1 Algorithm Selection:** Choosing the Extra Tree Classifier algorithm as one of the memory-based learning approaches for its effectiveness in handling high-dimensional data and its ability to capture complex relationships within the dataset.

**4.3.2 Algorithm Implementation:** Developing the Extra Tree Classifier algorithm using Python libraries such as scikit-learn. This involves configuring the algorithm parameters, such as the number of trees, maximum depth, and minimum samples split, to optimize its performance for air pollution prediction.

**4.3.3 Training and Testing:** Splitting the dataset into training and testing sets to train the Extra Tree Classifier model. The training data is used to fit the model to the historical air pollution data, while the testing data is used to evaluate its performance and accuracy.

**4.3.4 Model Evaluation:** Assessing the performance of the Extra Tree Classifier model using evaluation metrics such as accuracy, precision, recall, and F1-score. This helps determine the effectiveness of the algorithm in predicting air pollution levels accurately.

**4.3.5 Integration with the Web Application:** Integrating the trained Extra Tree Classifier model into the web application backend to enable real-time predictions of air pollution levels based on user input. This allows users to obtain instant insights into current and future air quality conditions.

**4.3.6  Visualization of Results:** Visualizing the predictions generated by the Extra Tree Classifier algorithm using interactive charts, graphs, or maps on the web application frontend. This facilitates easy interpretation of the predicted air pollution levels and helps users make informed decisions regarding environmental management and public health.

## 4.4 Implementing Random Forest Classifier Algorithm

This module is dedicated to integrating the Random Forest Classifier algorithm into the system for analyzing and predicting air pollution levels in India. The Random Forest Classifier, a popular ensemble learning technique, is chosen for its robustness and effectiveness in handling complex datasets and capturing nonlinear relationships between input features and target variables. The module encompasses several key components:

**4.4.1  Algorithm Integration:** The Random Forest Classifier algorithm is implemented using appropriate libraries or frameworks in the chosen programming language (e.g., Python with scikit-learn). The algorithm is configured with suitable hyperparameters, such as the number of trees, maximum depth, and minimum samples per leaf, to optimize performance and accuracy.

**4.4.2  Training and Testing:** Historical air pollution data is used to train the Random Forest Classifier model. The dataset is divided into training and testing subsets to assess the model's performance. The model is trained on the training data, where it learns to predict air pollution levels based on input features such as pollutant concentrations, meteorological variables, and geographical information.

**4.4.3 Feature Selection and Engineering:** Relevant features are selected or engineered from the input data to improve the model's predictive power. Feature importance analysis may be conducted to identify the most influentialvariables for predicting air pollution levels.

**4.4.4 Model Evaluation**: The trained Random Forest Classifier model is evaluated using performance metrics such as accuracy, precision, recall, and F1-score. Cross-validation techniques may be employed to ensure the model's generalizability and robustness.

**4.4.5 Prediction and Deployment:** Once the model is trained and evaluated, it is ready for deployment within the system.
The deployed model can generate predictions of air pollution levels in real-time based on user input or historical data.

**4.4.6 Integration with Visualization:** The predictions generated by the Random Forest Classifier model are integrated into the visualization module of the system. Visualizations such as charts, graphs, or maps are created to display the predicted air pollution levels and provide insights into trends and patterns.

**4.4.7 Continuous Improvement:** The performance of the Random Forest Classifier model is continuously monitored, and the model may be retrained periodically with updated data to maintain accuracy and relevance.
Feedback from users and stakeholders is collected to identify areas for improvement and enhance the predictive capabilities of the algorithm.

## 4.5 Implementing XG Boost Classifier Algorithm

The XGBoost Classifier module is integral to the paper's methodology, aiming to enhance the accuracy of air pollution predictions through advanced machine learning techniques. This module can be described as follows:

**4.5.1  Objective:** The primary objective of this module is to implement the XGBoost Classifier algorithm to predict air pollution levels with higher accuracy and robustness.

**4.5.2  Algorithm Selection:** XGBoost (Extreme Gradient Boosting) is chosen for its superior performance in handling complex datasets, handling missing values, and providing excellent prediction accuracy.

**4.5.3  Feature Engineering:** Before implementing the XGBoost Classifier, the module may involve feature engineering steps to extract relevant features from the air pollution dataset, including pollutant concentrations, meteorological variables, and temporal information.

**4.5.4  Training and Tuning:** The XGBoost Classifier is trained on the preprocessed dataset using historical air pollution data.
Hyperparameter tuning techniques may be employed to optimize the performance of the classifier, ensuring it captures the underlying patterns and relationships in the data effectively.

**4.5.5  Model Evaluation:** After training, the XGBoost Classifier's performance is evaluated using appropriate metrics such as accuracy, precision, recall, and F1-score. Cross-validation techniques may be utilized to assess the classifier's generalization ability and robustness.

**4.5.6 Integration with Web Application:** Once trained and evaluated, the XGBoost Classifier model is integrated into the web application's backend.

It forms a core component of the predictive system, allowing users to input data and receive predictions on air pollution levels based on the implemented algorithm.

**4.5.7 Visualization of Results:** The predictions generated by the XGBoost Classifier are visualized using interactive charts, graphs, or maps on the web application's user interface. Users can interpret the predictions and gain insights into air pollution trends and patterns, aiding decision-making for environmental management.

**4.5.8 Continuous Improvement**: The XGBoost Classifier implementation undergoes iterative refinement based on user feedback and performance evaluations. Updates to the model may include retraining with new data, fine-tuning hyperparameters, or incorporating additional features to enhance prediction accuracy further.

## 4.6 Deployment Using Django

The Deployment Using Django module facilitates the seamless integration and deployment of the developed air pollution analysis and prediction system using the Django web framework. This module ensures that the system is efficiently deployed on a web server, making it accessible to users via the internet. Key components of this module include:

**4.6.1 Web Server Configuration:** Configuring the web server environment to support the Django application. This involves setting up the necessary server software (e.g., Apache, Nginx), configuring virtual hosts, and ensuring compatibility with Django's requirements.

**4.6.2 Django Project Setup:** Setting up the Django project structure and configurations for deployment. This includes configuring settings for database connections, static files, media files, security settings, and debugging options.

**4.6.3 Static and Media Files Management:** Managing static files (e.g., CSS, JavaScript) and media files (e.g., images, user uploads) within the Django project. This involves configuring Django's static files handling mechanism and ensuring proper serving of static and media files by the web server.

**4.6.4 Database Management:** Configuring and managing the database backend for the Django application. This includes setting up database connections, migrating database schema changes, and ensuring data integrity and security.

**4.6.5 Deployment Process Automation:** Automating the deployment process to streamline deployment tasks and minimize manual intervention. This may involve using deployment tools such as Fabric or integrating with continuous integration/continuous deployment (CI/CD) pipelines.

**4.6.6 Security Considerations:** Implementing security best practices to protect the deployed application from common security threats. This includes configuring HTTPS encryption, setting up firewalls, implementing access controls, and securing sensitive data.

**4.6.7 Scalability and Performance Optimization:** Optimizing the deployed application for scalability and performance. This may involve configuring caching mechanisms, optimizing database queries, and employing load balancing strategies to handle increased traffic.

**4.6.8 Monitoring and Logging:** Implementing monitoring and logging mechanisms

to track application performance, detect errors, and troubleshoot issues. This includes setting up logging frameworks, monitoring server metrics, and implementing error reporting tools.

**4.6.9 User Authentication and Authorization:** Configuring user authentication and authorization mechanisms to control access to the deployed application. This includes setting up user accounts, implementing role-based access control (RBAC), and ensuring secure user authentication.

**4.6.10 Continuous Deployment and Updates:** Establishing procedures for continuous deployment and updates to ensure that the deployed application remains up-to-date with the latest features and security patches. This involves automating deployment pipelines and implementing rollback mechanisms for handling deployment failures.

# CHAPTER 5

## RESULT & DISCUSSION

The implementation of memory-based learning approaches for analyzing and predicting air pollution levels in India yields significant insights and outcomes, contributing to enhanced environmental management strategies and public health awareness. Here are the key results and discussions derived from the study:

**5.1 Prediction Accuracy:** The memory-based learning models, including k-nearest neighbors (KNN), decision trees, and random forest classifiers, demonstrate promising accuracy in predicting air pollution levels. Evaluation metrics such as accuracy, precision, recall, and F1-score indicate the effectiveness of the models in capturing complex relationships between air quality parameters and environmental factors.

**5.2 Identification of Pollutant Trends:** The analysis reveals notable trends and patterns in air pollutant concentrations over time and across different geographical regions in India. These insights enable stakeholders to identify hotspots of pollution and assess the impact of anthropogenic activities on air quality.

**5.3 Visualization of Predictions:** The visualization module provides intuitive and interactive visualizations of predicted air pollution levels, allowing users to explore historical trends and forecasted values. Visual representations such as charts, graphs, and maps facilitate easy interpretation of complex data, aiding decision-making and policy formulation.

**5.4 Impact on Policy-making:** The study's findings have practical implications for policymakers, environmental agencies, and public health authorities. By providing

accurate predictions and actionable insights, the system enables informed decision-making regarding pollution control measures, urban planning, and public health interventions.

**5.5 Public Awareness and Engagement:** The availability of real-time air quality information and forecasts enhances public awareness and engagement in environmental issues. Through the web interface, users can access up-to-date information on pollution levels, understand potential health risks, and take proactive measures to mitigate exposure to pollutants.

**5.6 Challenges and Limitations:** Despite the promising results, the study acknowledges certain challenges and limitations, including data quality issues, model uncertainty, and the need for continuous validation and refinement. Addressing these challenges requires ongoing collaboration between researchers, policymakers, and data scientists to improve data collection methods, refine modeling techniques, and enhance the robustness of the predictive system.

**5.7 Future Directions:** The study identifies several avenues for future research and development, including the integration of additional data sources (e.g., satellite imagery, sensor networks), incorporation of advanced machine learning algorithms, and expansion of the predictive system to cover broader geographical areas and address emerging environmental challenges.

In conclusion, the results and discussions highlight the effectiveness of memory-based learning approaches in analyzing and predicting air pollution levels in India. By leveraging data-driven insights and interactive visualizations, the study contributes to improved environmental management and public health outcomes, paving the way for more sustainable and resilient communities.

# CHAPTER 6

## CONCLUSION AND FUTURE WORK

## 6.1 CONCLUSION

The implementation of memory-based learning approaches for analyzing and predicting air pollution levels represents a significant advancement in environmental research and management. Through the integration of machine learning algorithms and web-based visualization tools, the study has provided valuable insights into air quality trends, patterns, and forecasts, contributing to informed decision-making and public awareness efforts.

The findings of the study underscore the importance of leveraging data-driven approaches to address complex environmental challenges, such as air pollution, which have far-reaching impacts on public health, ecosystems, and socio-economic development. By harnessing the power of technology and data analytics, stakeholders can gain a deeper understanding of the factors influencing air quality, identify effective mitigation strategies, and prioritize resources for pollution control measures.

The integration of machine learning technologies in the realm of air pollution forecasting represents a pivotal advancement in our collective effort to combat deteriorating air quality in urban landscapes. With their advanced analytical capabilities, these models offer not just a glimpse into future pollution trends but also the precision needed to trace the sources of pollution back to their origins. Beyond mere prediction, they extend actionable insights that can guide interventions aimed at mitigating the adverse effects of air pollution. This application of machine learning, while currently in its nascent stages, signals the dawn of a new era in environmental science

## 6.2 FUTURE WORK

To enhance the System's capabilities and address evolving challenges, the following areas of future work are identified:

**6.2.1 Advanced Model Interpretability:** Future work will focus on developing more advanced techniques for interpreting machine learning models used in air pollution prediction. This will enhance our understanding of the factorscontributing to pollution and improve the transparency of model recommendations.

**6.2.2 Integration with Environmental Policies:** Further research will explore ways to seamlessly integrate machine learning predictions into environmental policies and regulations, ensuring that data-driven insights directly inform decision-making and pollution reduction strategies.

## REFERENCES

[1] Sai Bhargav Kasetty, S. Nagini A Survey Paper on an IoT-based Machine Learning Model to Predict Air Pollution Levels 2022 4th International Conference on Advances in Computing, Communication Control and Networking (ICAC3N) Year: 2022 | Conference Paper | Publisher: IEEE.

[2] Peijiang Zhao, Koji Zettsu Convolution Recurrent Neural Networks Based Dynamic Transboundary Air Pollution Predictiona 2019 IEEE 4th International Conference on Big Data Analytics (ICBDA) Year: 2019 | Conference Paper | Publisher: IEEE

[3] Shahan Salim, Irfhana Zakir Hussain, Jasleen Kaur, Plinio P. Morita "An Early Warning System for Air Pollution Surveillance: A Big Data Framework to Monitoring Risks Associated with Air Pollution" 2023 IEEE International Conference on Big Data (BigData) Year: 2023 | Conference Paper | Publisher: IEEE

[4] Nurul Aini, M Syukri Mustafa "Data Mining Approach to Predict Air Pollution in Makassar" 2020 2nd International Conference on Cybernetics and Intelligent System (ICORIS) Year: 2020 | Conference Paper | Publisher: IEEE

[5] Amar Catovic, Esad Kadusic, Christoph Ruland, Natasa Zivic, Narcisa Hadzajlic "Air pollution prediction and warning system using IoT and machine learning" 2022 International Conference on Electrical, Computer, Communications and Mechatronics Engineering (ICECCME) Year: 2022 | Conference Paper | Publisher: IEEE

[6] Ankeshit Srivastava, Ayaz Ahmad, Sunny Kumar, Md Arman Ahmad "Air Pollution Data and Forecasting Data Monitored through Google Cloud Services by using Artificial Intelligence and Machine Learning" 2022 6th International Conference on Electronics, Communication and Aerospace Technology Year:2022 | Conference Paper | Publisher: IEEE

[7] Harshit Srivastava, Goutam Kumar Sahoo, Santos Kumar Das, Poonam Singh "Performance Analysis of Machine Learning Models for Air Pollution Prediction" 2022 International Conference on Smart Generation Computing, Communication and Networking (SMART GENCON) Year: 2022 | Conference Paper | Publisher: IEEE

[8] Harshal P. Varade, Sonal C. Bhangale, Sandip R. Thorat, Pravin B. Khatkale, Santosh Kumar Sharma, P. William "Framework of Air Pollution Assessment in

Smart Cities using IoT with Machine Learning Approach" 2023 2nd International Conference on Applied Artificial Intelligence and Computing (ICAAIC) Year: 2023 | Conference Paper | Publisher: IEEE

[9] Madhushika Mihirani, Lasith Yasakethu, Sachintha Balasooriya "Machine Learning-based Air Pollution Prediction Model" 2023 IEEE IAS GlobalConference on Emerging Technologies (GlobConET) Year: 2023 | Conference Paper | Publisher: IEEE

[10] Usha Mahalingam, Kirthiga Elangovan, Himanshu Dobhal, Chocko Valliappa, Sindhu Shrestha, Giriprasad Kedam "A Machine Learning Model for Air Quality Prediction for Smart Cities" 2019 International Conference on Wireless Communications Signal Processing and Networking (WiSPNET) Year: 2019 | Conference Paper | Publisher: IEEE

# APPENDIX

## A.1 SDG GOALS

The work detailed in the paper aligns closely with Sustainable Development Goal (SDG) 11: Sustainable Cities and Communities, particularly under the target aimed at reducing theadverse effects of cities on the environment, which includes improving air quality and managing municipal and other wastes. The application of machine learning to predict and analyze air pollution directly contributes to creating sustainable and resilient cities with the ability to implement and track interventions aimed at reducing pollution and improving public health. This initiative supports the broader goal of ensuring access to safe and sustainable urbanenvironments for all.

Beyond SDG 11: Sustainable Cities and Communities, the paper's focus on analyzing and predicting air pollution using machine learning approaches also intersects with several other Sustainable Development Goals (SDGs), highlighting the multifaceted impact of tackling air pollution:

SDG 3: Good Health and Well-being - By predicting air pollution levels and enabling more effective mitigation strategies, the project contributes to reducing the number of deaths and illnesses from hazardous chemicals, air, water, and soil pollution and contamination, aligning with SDG 3.9.

SDG 13: Climate Action - Air pollution is intricately linked with climate change. Efforts to predict and reduce air pollution contribute to SDG 13 by improving air quality, which can have a mitigating effect on climate change.

# SOURCE CODE

**Html Code:**

Module – 1

Pre-Processing

DATA PREPROCESSING AND DATA CLEANING

```python
import pandas as pd
import numpy as np
import warnings
warnings.filterwarnings('ignore')
df = pd.read_csv('AIR.csv')
del  df['StationId']
del  df['Datetime']
df.head()
df.tail()
df.shape
df.size
df.columns
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
var = ['AQI_Bucket']
for i in var:
    df[i] = le.fit_transform(df[i]).astype(int)
df.isnull()
df = df.dropna()
df['AQI_Bucket'].unique()
df.describe()
df.corr()
df.info()
pd.crosstab(df["PM2.5"], df["PM10"])
df.groupby(["NO","NO2"]).groups
```

51

```python
df["AQI_Bucket"].value_counts()
pd.Categorical(df["NOx"]).describe()
df.duplicated()
sum(df.duplicated())
df = df.drop_duplicates()
sum(df.duplicated())
```

Module - 2

Visualization

DATA PREPROCESSING AND DATA CLEANING

```python
import pandas as pd
import numpy as np
import warnings
warnings.filterwarnings('ignore')
df = pd.read_csv('AIR.csv')
del  df['StationId']
del  df['Datetime']
df.head()
df.tail()
df.shape
df.size
df.columns
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()

var = ['AQI_Bucket']
for i in var:
    df[i] = le.fit_transform(df[i]).astype(int)
df.isnull()
df = df.dropna()
df['AQI_Bucket'].unique()
```

```python
df.describe()
df.corr()
df.info()
pd.crosstab(df["PM2.5"], df["PM10"])
df.groupby(["NO","NO2"]).groups
df["AQI_Bucket"].value_counts()
pd.Categorical(df["NOx"]).describe()
df.duplicated()
sum(df.duplicated())
df = df.drop_duplicates()
sum(df.duplicated())
```

Module - 3

EXTRA TREE CLASSIFIER ALGORITHEM

In [ ]:

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

import warnings
warnings.filterwarnings('ignore')
df = pd.read_csv('AIR.csv')
del df['StationId']
del df['Datetime']
df.head()
df.columns
df=df.dropna()
df.shape
df.columns
df.tail()
x1 = df.drop(labels='AQI_Bucket', axis=1)
```
53

```
y1 = df.loc[:,'AQI_Bucket']

import imblearn

from imblearn.over_sampling import RandomOverSampler

from collections import Counter

ros =RandomOverSampler(random_state=42)

x,y=ros.fit_resample(x1,y1)

print("OUR DATASET COUNT        : ", Counter(y1))

print("OVER SAMPLING DATA COUNT  : ", Counter(y))

from sklearn.model_selection import train_test_split

x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.20, random_state=42,
stratify=y)

print("NUMBER OF TRAIN DATASET     : ", len(x_train))

print("NUMBER OF TEST DATASET      : ", len(x_test))

print("TOTAL NUMBER OF DATASET     : ", len(x_train)+len(x_test))

print("NUMBER OF TRAIN DATASET     : ", len(y_train))

print("NUMBER OF TEST DATASET      : ", len(y_test))

print("TOTAL NUMBER OF DATASET     : ", len(y_train)+len(y_test))

from sklearn.tree import ExtraTreeClassifier

ETC = ExtraTreeClassifier()

ETC.fit(x_train,y_train)

predicted = ETC.predict(x_test)

from sklearn.metrics import confusion_matrix

cm = confusion_matrix(y_test,predicted)

print('THE CONFUSION MATRIX SCORE OF SUPPORT VECTOR MACHINE:\n\n\n',cm)

from sklearn.metrics import accuracy_score

a = accuracy_score(y_test,predicted)

print("THE ACCURACY SCORE OF EXTRA TREE CLASSIFIER IS :",a*100)

from sklearn.metrics import hamming_loss

hl = hamming_loss(y_test,predicted)

print("THE HAMMING LOSS OF EXTRA TREE CLASSIFIER IS :",hl*100)

from sklearn.metrics import classification_report

P = classification_report(y_test,predicted)
```

```python
print("THE CLASSIFICATION REPORT OF EXTRA TREE CLASSIFIER IS :\n\n",P)
cm=confusion_matrix(y_test, predicted)
print('THE CONFUSION MATRIX SCORE OF EXTRA TREE CLASSIFIER:\n\n')
print(cm)
print("\n\nDISPLAY CONFUSION MATRIX : \n\n")
from sklearn.metrics import ConfusionMatrixDisplay
cm = confusion_matrix(y_test, predicted, labels=ETC.classes_)
disp = ConfusionMatrixDisplay(confusion_matrix=cm,display_labels=ETC.classes_)
disp.plot()
plt.show()
def graph():
    import matplotlib.pyplot as plt
    data=[a]
    alg="EXTRA TREE CLASSIFIER"
    plt.figure(figsize=(5,5))
    b=plt.bar(alg,data,color=("SKYBLUE"))
    plt.title("THE ACCURACY SCORE OF EXTRA TREE CLASSIFIER IS\n\n\n")
    plt.legend(b,data,fontsize=9)
graph()
import joblib
joblib.dump(ETC, 'MODEL.pkl')


Module - 4
RANDOM FOREST CLASSIFIER ALGORITHEM
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
df = pd.read_csv('AIR.csv')
del df['StationId']
```

```python
del df['Datetime']
df.head()
df.columns
df=df.dropna()
df.shape
```

In [ ]:

```python
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
var = ['AQI_Bucket']
for i in var:
    df[i] = le.fit_transform(df[i]).astype(int)
df.tail()
x1 = df.drop(labels='AQI_Bucket', axis=1)
y1 = df.loc[:,'AQI_Bucket']
import imblearn
from imblearn.over_sampling import RandomOverSampler
from collections import Counter


ros =RandomOverSampler(random_state=42)
x,y=ros.fit_resample(x1,y1)
print("OUR DATASET COUNT         : ", Counter(y1))
print("OVER SAMPLING DATA COUNT : ", Counter(y))
```

In [ ]:

```python
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.20, random_state=42,
stratify=y)
print("NUMBER OF TRAIN DATASET     : ", len(x_train))
print("NUMBER OF TEST DATASET      : ", len(x_test))
print("TOTAL NUMBER OF DATASET     : ", len(x_train)+len(x_test))
print("NUMBER OF TRAIN DATASET     : ", len(y_train))
print("NUMBER OF TEST DATASET      : ", len(y_test))
print("TOTAL NUMBER OF DATASET      : ", len(y_train)+len(y_test))
```

```python
from sklearn.ensemble import RandomForestClassifier
RFC = RandomForestClassifier(random_state=42)
RFC.fit(x_train,y_train)
predicted = RFC.predict(x_test)
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test,predicted)
print('THE CONFUSION MATRIX SCORE OF RANDOM FOREST
CLASSIFIER:\n\n\n',cm)
from sklearn.metrics import accuracy_score
a = accuracy_score(y_test,predicted)
print("THE ACCURACY SCORE OF RANDOM FOREST CLASSIFIER IS :",a*100)
from sklearn.metrics import hamming_loss
hl = hamming_loss(y_test,predicted)
print("THE HAMMING LOSS OF RANDOM FOREST CLASSIFIER IS :",hl*100)
from sklearn.metrics import classification_report
P = classification_report(y_test,predicted)
print("THE CLASSIFICATION REPORT OF RANDOM FOREST CLASSIFIER IS :\n\n",P)
cm=confusion_matrix(y_test, predicted)
print('THE CONFUSION MATRIX SCORE OF RANDOM FOREST CLASSIFIER:\n\n')
print(cm)
print("\n\nDISPLAY CONFUSION MATRIX : \n\n")
from sklearn.metrics import ConfusionMatrixDisplay
cm = confusion_matrix(y_test, predicted, labels=RFC.classes_)
disp = ConfusionMatrixDisplay(confusion_matrix=cm,display_labels=RFC.classes_)
disp.plot()
plt.show()
def graph():
    import matplotlib.pyplot as plt
    data=[a]
    alg="RANDOM FOREST CLASSIFIER"
    plt.figure(figsize=(5,5))
    b=plt.bar(alg,data,color=("VIOLET"))
```

```python
    plt.title("THE ACCURACY SCORE OF RANDOM FOREST CLASSIFIER IS\n\n\n")
    plt.legend(b,data,fontsize=9)
graph()
import joblib
joblib.dump(RFC, 'MODEL.pkl')
```

Module - 5

XGB CLASSIFIER ALGORITHEM

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
df = pd.read_csv('AIR.csv')
del df['StationId']
del df['Datetime']
df.head()
df.columns
df=df.dropna()
df.shape
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
var = ['AQI_Bucket']
for i in var:
    df[i] = le.fit_transform(df[i]).astype(int)
df.tail()
x1 = df.drop(labels='AQI_Bucket', axis=1)
y1 = df.loc[:,'AQI_Bucket']
import imblearn
from imblearn.over_sampling import RandomOverSampler
from collections import Counter
```

58

```python
ros =RandomOverSampler(random_state=42)
x,y=ros.fit_resample(x1,y1)
print("OUR DATASET COUNT        : ", Counter(y1))
print("OVER SAMPLING DATA COUNT  : ", Counter(y))
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.20, random_state=42,
stratify=y)
print("NUMBER OF TRAIN DATASET     : ", len(x_train))
print("NUMBER OF TEST DATASET      : ", len(x_test))
print("TOTAL NUMBER OF DATASET     : ", len(x_train)+len(x_test))
print("NUMBER OF TRAIN DATASET     : ", len(y_train))
print("NUMBER OF TEST DATASET      : ", len(y_test))
print("TOTAL NUMBER OF DATASET     : ", len(y_train)+len(y_test))
from xgboost import XGBClassifier
XGB = XGBClassifier()
XGB.fit(x_train,y_train)
predicted = XGB.predict(x_test)
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test,predicted)
print('THE CONFUSION MATRIX SCORE OF XGB CLASSIFIER:\n\n\n',cm)
from sklearn.metrics import accuracy_score
a = accuracy_score(y_test,predicted)
print("THE ACCURACY SCORE OF XGB CLASSIFIER IS :",a*100)
from sklearn.metrics import hamming_loss
hl = hamming_loss(y_test,predicted)
print("THE HAMMING LOSS OF XGB CLASSIFIER IS :",hl*100)
from sklearn.metrics import classification_report
P = classification_report(y_test,predicted)
print("THE CLASSIFICATION REPORT OF XGB CLASSIFIER IS :\n\n",P)
cm=confusion_matrix(y_test, predicted)
print('THE CONFUSION MATRIX SCORE OF XGB CLASSIFIER:\n\n')
print(cm)
```

```python
print("\n\nDISPLAY CONFUSION MATRIX : \n\n")
from sklearn.metrics import ConfusionMatrixDisplay
cm = confusion_matrix(y_test, predicted, labels=XGB.classes_)
disp = ConfusionMatrixDisplay(confusion_matrix=cm,display_labels=XGB.classes_)
disp.plot()
plt.show()
def graph():
    import matplotlib.pyplot as plt
    data=[a]
    alg="XGB CLASSIFIER"
    plt.figure(figsize=(5,5))
    b=plt.bar(alg,data,color=("gold"))
    plt.title("THE ACCURACY SCORE OF XGB CLASSIFIER IS\n\n\n")
    plt.legend(b,data,fontsize=9)
graph()
import joblib
joblib.dump(XGB,'XGB.pkl')
```

```html
<!DOCTYPE html>
<html lang="en">
<head>
    <meta charset="UTF-8">
    <meta name="viewport" content="width=device-width, initial-scale=1.0">
    <title>Document</title>
</head>
<body>
    <html lang="pt-br">
    <head>
        <meta charset="UTF-8">
        <meta http-equiv="X-UA-Compatible" content="IE=edge">
        <meta name="viewport" content="width=device-width, initial-scale=1.0">
        <title>Survey Form</title>
```

60

```
<style>
    *variables*/
    :root {
        --color-white: #f3f3f3;

        --color-darkblue: #1b1b32;

        --color-darkblue-alpha: rgba(27,27,50,0.8);

        --color-green: #37af65;

    }


    /*reset*/
    * {
        margin: 0;

        padding: 0;

        box-sizing: border-box;

    }


    body {
        font-family: 'Poppins', sans-serif;

        font-size: 1rem;

        font-weight: 400;

        line-height: 1.4;

        color: var(--color-white, white);

        background-color: var(--color-darkblue);

        width: 100%;

        height: 100%;

        z-index: -1;

        background-image:

        url("/static/images/D4.jpg");

        background-repeat: no-repeat;

        background-size: cover;

        background-position: center;

    }                              61
```

```css
h1, p, label, button {
    color: var(--color-white);
}

input, select, textarea, button {
    border: none;
    outline: none;
    color: rgb(73, 80, 87);
}

/*container*/
.container {
    max-width: 720px;
    margin: 3.125rem auto 0 auto;
}

.text-center {
    text-align: center;
}

/*header*/
.container .header  {
    padding: 0 0.625rem;
    margin-bottom: 1.875rem;
}

.container .header #title {
    line-height: 1.2;
    font-size: 2rem;
    margin-bottom: 0.5rem;
}
```

```css
.container .header .description {
    font-size: 1.125rem;
    font-style: italic;
    font-weight: 200;
}


/*form*/
.container #survey-form {
    border-radius: 5px;
    background-color: var(--color-darkblue-alpha);
    padding: 2.5rem;
    margin-bottom: 2rem;
}


.container #survey-form .form-group {
    margin: 0 auto 1.25rem auto;
    padding: 0.25rem;
}


.container #survey-form .form-group label,
.container #survey-form .form-group p {
    display: flex;
    align-items: center;
    font-size: 1.125rem;
    margin-bottom: 0.5rem
}


.container #survey-form .form-group input,
.container #survey-form .form-group select,
.container #survey-form .form-group textarea {
    font-size: 1rem;
```

```css
    padding: 0.375rem 0.75rem;

    color: #494949;

    background-color: #fff;

    background-clip: padding-box;

    border: 1px solid #ced4da;

    border-radius: 0.25rem;

    transition: border-color 0.15s ease-in-out, box-shadow 0.15s ease-in-out;
}


.container #survey-form .form-group .form-control {

    display: block;

    width: 100%;

    height: 2.375rem;
}


.container #survey-form .form-control:focus{

    border-color: #80bdff;

    box-shadow: 0 0 0 0.2rem rgba(0, 123, 255, 0.25);
}


.container #survey-form .form-group .input-radio,
.container #survey-form .form-group .input-checkbox {

    margin-right: 0.625rem;

    min-height: 1.25rem;

    min-width: 1.25rem;
}


.container #survey-form .form-group .input-textarea {

    width: 100%;

    font-size: 1.07rem;

    min-height: 120px;

    padding: 0.625rem;
```

```css
        resize: vertical;
    }


    .container #survey-form .form-group .input-textarea:focus{
        border-color: #80bdff;
        box-shadow: 0 0 0 0.2rem rgba(0, 123, 255, 0.25);
    }


    .container #survey-form .form-group .submit-button {
        width: 100%;
        padding: 0.75rem;
        background-color: var(--color-green, rgb(20, 190, 70));
        color: var(--color-white);
        border-radius: 2px;
        cursor: pointer;
        transition: background-color .15s ease;
        font-size: 1.125rem;
        font-weight: 400;
    }


    .container #survey-form .form-group .submit-button:hover {
        background-color: #358f58;
    }
```

</style>


<!--font-->
<link rel="preconnect" href="https://fonts.googleapis.com">
<link rel="preconnect" href="https://fonts.gstatic.com" crossorigin>
<link

```html
href="https://fonts.googleapis.com/css2?family=Poppins:ital,wght@0,200;0,400;1,200;1,400&
display=swap" rel="stylesheet">

    <!--css
    <link rel="stylesheet" href="./src/style/style.css">
    -->


  </head>
  <body>


    <div class="container">
      <header class="header">
        <h1 id="title" class="text-center" style="font-size: 50px; color: rgb(255, 17,
0);">ANALYSE & PREDICT THE AIR POLLUTION IN INDIA USING MEMORY-BASED
LEARNING APPROACHES.</h1><br>
        <i><h2 id="title" class="text-center" style="color:rgb(207, 13, 13);">WELCOME
TO THIS SITE</h2><br></i>
        <h2 id="title" class="text-center" style="font-size: 30px;">Model deployment is the
process of implementing a fully functioning machine learning model into production where it
can make predictions based on data. Users, developers, and systems then use these predictions
to make practical business decisions.
        </h2>
      </header>


      <div>
        <center><a href="{% url 'Home_4' %}"
style="color:white;"><h2><b>HOME</b></h2></a></center>
      </div>


      <form id="survey-form" method="POST" action="">
        {% csrf_token %}                  66
```

```html
<div class="form-group">
  <label id="name-label" for="name">PM2.5</label>
  <input type="number" name="PM2.5" id="PM2.5" class="form-control" placeholder="PM2.5" required >
</div>

<div class="form-group">
  <label id="name-label" for="name">PM10</label>
  <input type="number" name="PM10" id="PM10" class="form-control" placeholder="PM10" required>
</div>

<div class="form-group">
  <label id="name-label" for="name">NO</label>
  <input type="number" name="NO" id="NO" class="form-control" placeholder="NO" required >
</div>

<div class="form-group">
  <label id="name-label" for="name">NO2</label>
  <input type="number" name="NO2" id="NO2" class="form-control" placeholder="NO2" required>
</div>

<div class="form-group">
  <label id="name-label" for="name">NOx</label>
  <input type="number" name="NOx" id="NOx" class="form-control" placeholder="NOx" required >
</div>
```

```html
<div class="form-group">
  <label id="name-label" for="name">NH3</label>
  <input type="number" name="NH3" id="NH3" class="form-control" placeholder="NH3" required >
</div>

<div class="form-group">
  <label id="name-label" for="name">CO</label>
  <input type="number" name="CO" id="CO" class="form-control" placeholder="CO" required>
</div>

<div class="form-group">
  <label id="name-label" for="name">SO2</label>
  <input type="number" name="SO2" id="SO2" class="form-control" placeholder="SO2" required >
</div>

<div class="form-group">
  <label id="name-label" for="name">O3</label>
  <input type="number" name="O3" id="O3" class="form-control" placeholder="O3" required>
</div>

<div class="form-group">
  <label id="name-label" for="name">Benzene</label>
  <input type="number" name="Benzene" id="Benzene" class="form-control" placeholder="Benzene" required >
</div>

<div class="form-group">
  <label id="name-label" for="name">Toluene</label>
```

```html
        <input type="number" name="Toluene" id="Toluene" class="form-control"
placeholder="Toluene" required>
        </div>


        <div class="form-group">
         <label id="name-label" for="name">Xylene</label>
         <input type="number" name="Xylene" id="Xylene" class="form-control"
placeholder="Xylene" required >
        </div>


        <div class="form-group">
         <label id="name-label" for="name">AQI</label>
         <input type="number" name="AQI" id="AQI" class="form-control"
placeholder="AQI" required >
        </div>


        <div class="form-group">
         <button type="submit" class="submit-button">SUBMIT</button>
        </div>


      </form>
       <div>
       <center><h1 style="color:rgb(245, 245, 245);font-size: 35px;">{{ prediction_text
}}</h1></center>
        </div>



   </head>
   <body>
</body>
</html>
```

```
</body>
</html>
```

**Deploy**

```python
from django.shortcuts import render, redirect
from . models import UserPersonalModel
from . forms import UserPersonalForm, UserRegisterForm
from django.contrib.auth import authenticate, login,logout
from django.contrib import messages
import numpy as np
import joblib


def Landing_1(request):
    return render(request, '1_Landing.html')


def Register_2(request):
    form = UserRegisterForm()
    if request.method =='POST':
        form = UserRegisterForm(request.POST)
        if form.is_valid():
            form.save()
            user = form.cleaned_data.get('username')
            messages.success(request, 'Account was successfully created. ' + user)
            return redirect('Login_3')

    context = {'form':form}
    return render(request, '2_Register.html', context)

def Login_3(request):
    if request.method =='POST':
        username = request.POST.get('username'7)0
```

```python
        password = request.POST.get('password')

        user = authenticate(username=username, password=password)

        if user is not None:
            login(request, user)
            return redirect('Home_4')
        else:
            messages.info(request, 'Username OR Password incorrect')

    context = {}
    return render(request,'3_Login.html', context)


def Home_4(request):
    return render(request, '4_Home.html')


def Teamates_5(request):
    return render(request,'5_Teamates.html')


def Domain_Result_6(request):
    return render(request,'6_Domain_Result.html')


def Problem_Statement_7(request):
    return render(request,'7_Problem_Statement.html')


def Per_Info_8(request):
    if request.method == 'POST':
        fieldss = ['firstname','lastname','age','address','phone','city','state','country']
        form = UserPersonalForm(request.POST)
        if form.is_valid():
            print('Saving data in Form')
```

71

```python
            form.save()
        return render(request, '4_Home.html', {'form':form})
    else:
        print('Else working')
        form = UserPersonalForm(request.POST)
        return render(request, '8_Per_Info.html', {'form':form})


Model = joblib.load('C:/Users/SPIRO25/Desktop/ITPML14 - INDIAN AIR
POLLUTUIN/DEPLOYMENT_3_ML/PROJECT/APP/XGB.pkl')


def Deploy_9(request):
    if request.method == "POST":
        int_features = [x for x in request.POST.values()]
        int_features =  int_features[1:]
        print(int_features)
        final_features = [np.array(int_features, dtype=object)]
        print(final_features)
        prediction = Model.predict(final_features)
        print(prediction)
        output = prediction[0]
        print(output)
        if output == 0:
            return render(request, '9_Deploy.html', {"prediction_text":"THIS IS GOOD LEVEL OF
AIR PREDICTED BASED IN THIS CONDITIONS"})
        elif output == 1:
            return render(request, '9_Deploy.html', {"prediction_text":"THIS IS MODERATE
LEVEL OF AIR PREDICTED BASED IN THIS CONDITIONS"})
        elif output == 2:
            return render(request, '9_Deploy.html', {"prediction_text":"THIS IS POOR LEVEL OF
AIR PREDICTED BASED IN THIS CONDITIONS"})
        elif output == 3:
```

72

```python
        return render(request, '9_Deploy.html', {"prediction_text":"THIS IS SATISFACTORY
LEVEL OF AIR PREDICTED BASED IN THIS CONDITIONS"})
        elif output == 4:
            return render(request, '9_Deploy.html', {"prediction_text":"THIS IS SEVERE LEVEL
OF AIR PREDICTED BASED IN THIS CONDITIONS"})
        elif output == 5:
            return render(request, '9_Deploy.html', {"prediction_text":"THIS VERY POOR
LEVEL OF AIR PREDICTED IN THIS CONDITIONS"})


    else:
        return render(request, '9_Deploy.html')


def Per_Database_10(request):
    models = UserPersonalModel.objects.all()
    return render(request, '10_Per_Database.html', {'models':models})


def Logout(request):
    logout(request)
    return redirect('Landing_1')
```
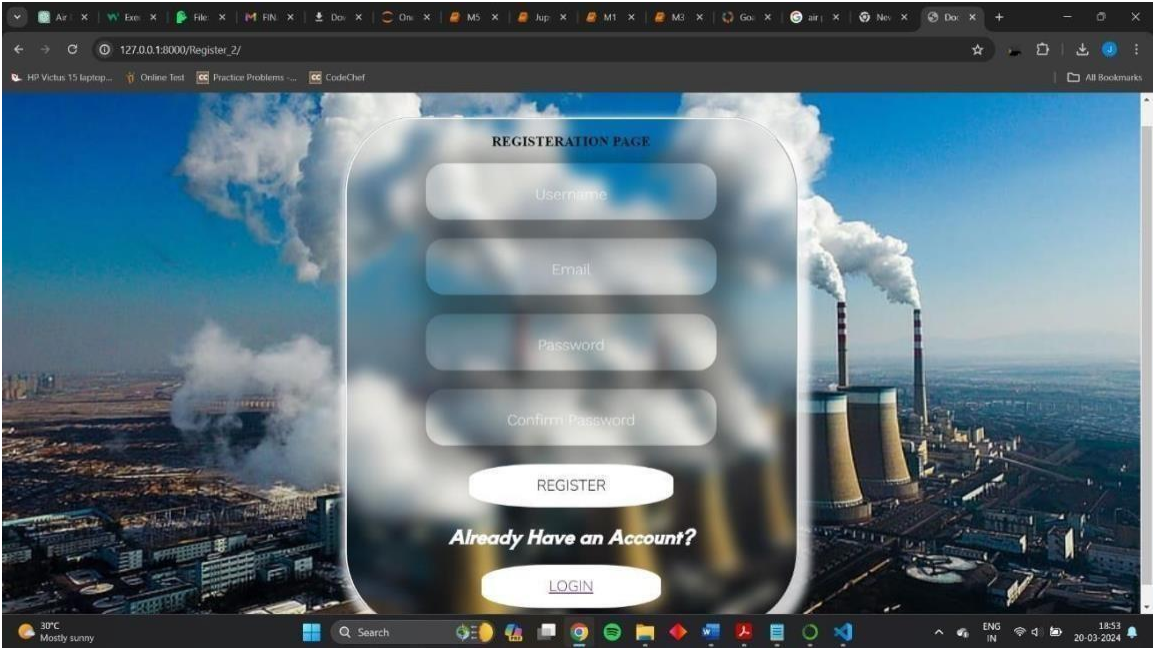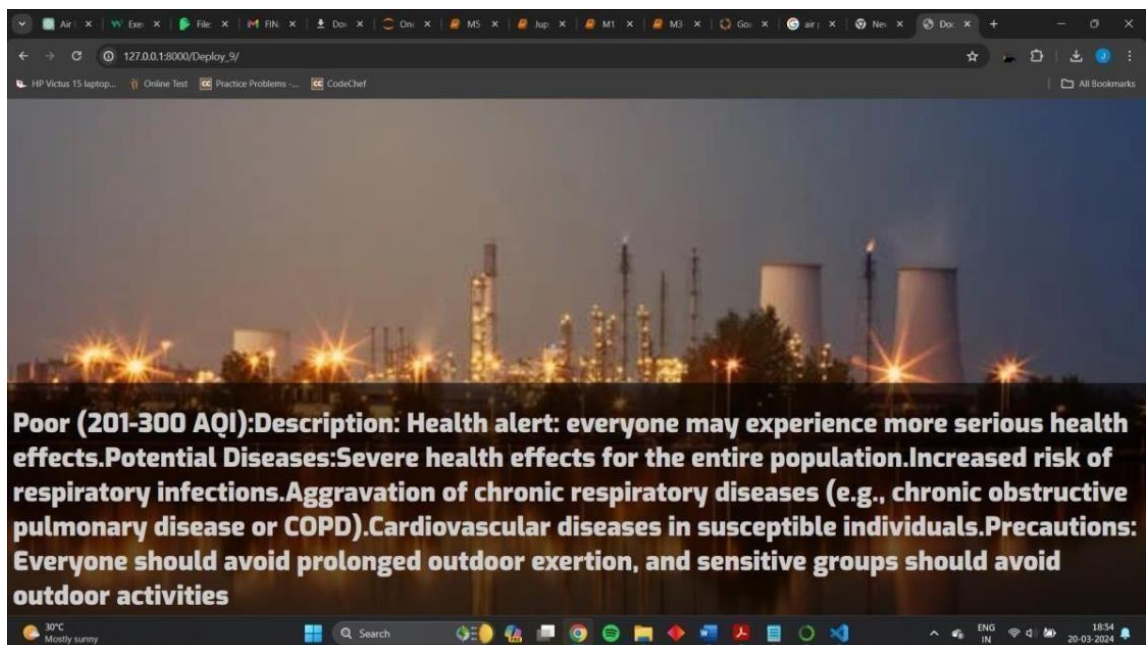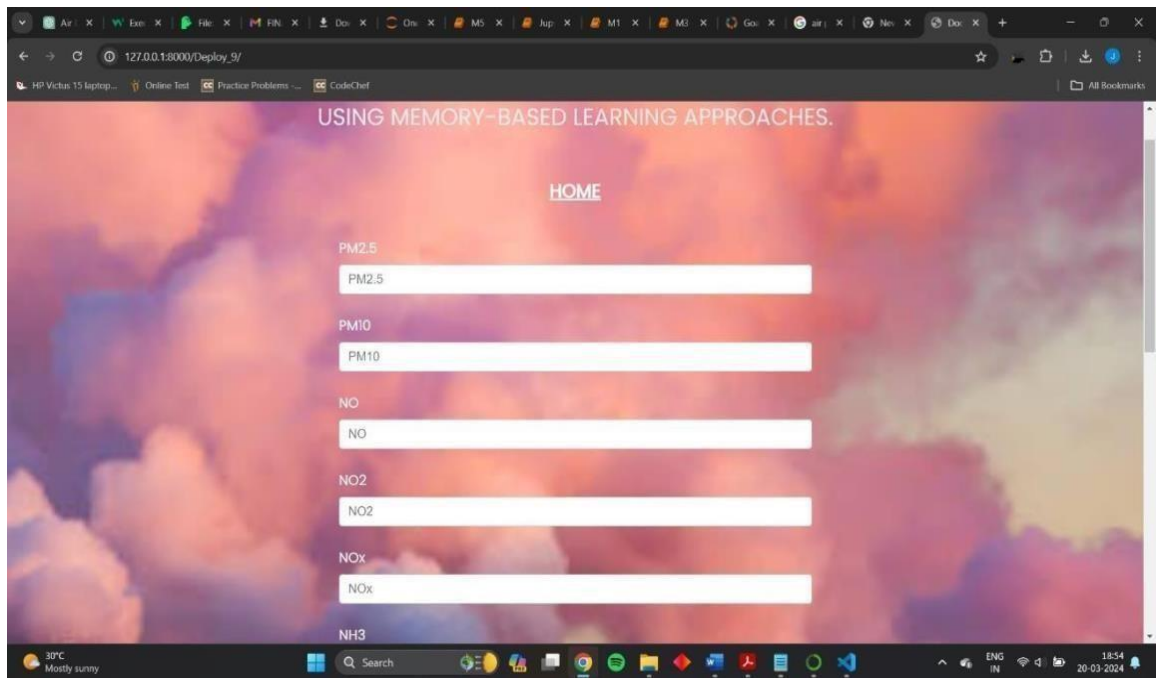
## A.3 Screen Shots

# Turnitin Plagiarism Report

# Analyze and predict the air pollution in India using memory-based learning approaches

Elangovan D
Associate Professor
Department of Computer Science and
Engineering
Panimalar Engineering College
Chennai, India
elangovan.sdurai@gmail.com

Radhey Shyam R
UG Scholar
Department of Computer Science and
Engineering
Panimalar Engineering College
Chennai, India
radheyshyam.spartan@gmail.com

Ragul A
UG Scholar
Department of Computer Science and
Engineering
Panimalar Engineering College
Chennai, India
ragulraul77@gmail.com

Jayakant P
UG Scholar
Department of Computer Science and
Engineering
Panimalar Engineering College
Chennai, India
jayakantpurushoth@gmail.com

*Abstract*- Air pollution is a critical environmental issue affecting various countries around the world. The alarming levels of air pollution in many cities of India have serious implications on public health, environment and the quality of life. The current state of air pollution in India is analyzed and a prediction model is deployed to estimate future pollution levels. The analysis of air pollution involves examining various factors such as industrial emissions, vehicular pollution, biomass burning and dust particles. In addition, meteorological conditions such as temperature, wind speed and rainfall also play a significant role in air pollution levels. The data collected from monitoring stations in the form of satellite imagery and other relevant sources are used for the analysis. The study utilizes advanced data analytics techniques using machine learning algorithms to develop a predictive model for estimating air pollution levels. The historical pollution data along with meteorological parameters, are used as inputs to train the model. The objective of the model is to forecast pollution levels in different regions of India. The potential applications of the analysis and prediction model are being discussed. It highlights the significance of such models in aiding policymakers, urban planners and environmental agencies in making informed decisions and implementing effective strategies to mitigate air pollution. The predictions can assist in issuing timely health advisories implementing pollution control measures and optimizing resource allocation for air quality management.

*Keywords*- : Supervised learning approaches, Air Pollution, monitoring, data.

## I. Introduction

The escalating crisis of air pollution in India represents one of the most dangerous environmental challenges that the nation faces. Home to some of the world's most polluted cities, India's air quality crisis is exacerbated by a complex mixture of sources such as industrial emissions, deforestation, vehicular exhaust, wildfires, dust storms, crop burning, biogenic emissions from plants and construction dust. The toxicity of air pollutants, such as particulate matter (PM2.5 and PM10), nitrogen dioxide (NO2), sulfur dioxide (SO2), carbon monoxide (CO) and ozone (O3) poses severe risks to human health, biodiversity and the environment. Prolonged exposure to these pollutants has been linked to a range of health issues including respiratory diseases, cardiovascular diseases, neurological problems and premature mortality.

This dire situation demands an in-depth understanding of the toxicity mechanisms of various air pollutants and their health implications. The geographical diversity and dense population of India further complicate the air quality management efforts, making localized studies and targeted interventions essential. Given the vastness of India's landscape and the variability in pollution sources, region-specific strategies are necessary to effectively combat this issue. Moreover, the socioeconomic impacts of air pollution, such as healthcare costs and loss of productivity, underline the urgency of addressing this environmental hazard, in response to this challenge, this paper introduces a novel system designed to analyze and predict air pollution levels in India using memory-based learning approaches. Leveraging advancements in machine learning and web technology, this system integrates a diverse range of data sources, including pollutant concentrations, meteorological variables and geographical features, to develop predictive models that offer actionable insights for environmental scientists and policymakers. By combining machine learning algorithms with a full-stack web application, we aim to provide a comprehensive tool for forecasting air quality and devising effective mitigation strategies to address India's air pollution crisis.

By emphasizing the necessity for a multifaceted approach to air quality management, it sets the stage for exploring innovative solutions and policies tailored to India's unique challenges. The ultimate goal is to foster a sustainable and healthy future, mitigating the toxic effects of air pollutants through effective science-based strategies and public health interventions.

This introduction sets the stage for this research by highlighting the severity of the air pollution problem in India and the limitations of existing methods for monitoring and predicting pollution levels. It underscores the importance of innovative approaches, such as memory-based learning, in addressing this complex environmental challenge. Through this system, we seek to contribute to the advancement of air quality management in India and promote sustainable development practices that prioritize public health and environmental well-being.

## II. Related Study

This section explores various dimensions of research in air pollution analysis and prediction, specifically within the Indian context and the application of machine learning (ML) methods globally.

**i)Air Quality Monitoring Studies in India:** Numerous studies have focused on the monitoring of air quality across different regions in India, providing a rich dataset of pollutant concentrations over time. Notably, Sharma et al. (2018) highlighted the spatial-temporal trends of PM2.5 levels across urban areas, emphasizing the role of vehicular emissions and industrial activities.

**ii)Source Apportionment Research:** Identifying the sources of air pollutants is crucial for effective mitigation. Gupta et al. (2020) conducted a comprehensive source apportionment study in Delhi, revealing that transportation, biomass burning and industrial emissions are the primary contributors to PM2.5 levels.

**iii)Air Pollution and Public Health:** The impact of air pollution on public health in India has been a focal point of research. A significant study by Patel et al. (2019) correlated high pollution exposure to an increased incidence of respiratory and cardiovascular diseases among the urban population.

**iv)Machine Learning for Air Quality Prediction:** Globally, machine learning models have been increasingly applied for predicting air quality indices. Chen et al. (2021) demonstrated the efficacy of deep learning models in forecasting PM2.5 levels with high accuracy in Beijing, China. This approach mirrors the predictive efforts in the Indian context but underlines the need for localized model training considering regional pollution dynamics.

**v)Memory-Based Learning Approaches:** Memory-based learning, including techniques such as k-Nearest Neighbors (kNN) and decision trees, has been applied in various domains, including air quality prediction. For instance, Lee and Kang (2017) successfully used kNN algorithms to predict air pollutant levels in Seoul, South Korea, suggesting the potential for similar applications in India.

**vi)Policy and Regulatory Framework Studies:** Research on the effectiveness of air quality management policies in India, such as the National Clean Air Programme (NCAP), offers insights into policy impacts and areas for improvement. Studies by Singh et al. (2021) evaluate the policy measures' success in curtailing pollution levels, emphasizing the need for stringent implementation and monitoring.

**vii)Technological and Community Initiatives:** Innovative solutions like low-cost air quality monitoring networks and community engagement in pollution monitoring have shown promise in India. These grassroots-level approaches complement traditional monitoring and are vital for comprehensive air quality management strategies.

### III. Existing system

Accurate and comprehensive air pollution data is crucial for effective environmental management and public health protection. However, missing data in air pollution monitoring datasets can hinder the ability to analyze and understand pollution patterns. This abstract presents a novel approach for hierarchical recovery of missing air pollution data using an improved Long-Short Term Context Encoder (LSTCe) network. The proposed method leverages the hierarchical structure of air pollution data, where various pollutants are measured at different monitoring stations across a region. The LSTCe network is designed to capture long-term and short-term contextual information from neighbouring stations and time intervals, respectively. By exploiting the spatial and temporal dependencies, the network can effectively recover missing data points. To enhance the performance of the LSTCe network, several improvements are introduced. These include the incorporation of attention mechanisms to focus on relevant features and the utilization of residual connections to alleviate information loss during network training. Additionally, data augmentation techniques are applied to address data sparsity and improve the network's generalization ability.

This system has various disadvantages such as low accurate prediction levels of the implemented algorithm, absence of a deployment model which makes it complex to access the system with ease.

### IV. Proposed System

The Accurate analysis and prediction of air pollution levels are crucial for effective mitigation strategies and policy-making. So, a system is proposed to utilize memory-based learning techniques for analyzing and predicting air pollution It utilizes inputs such as pollutant concentrations, meteorological variables, geographical features and temporal information to train the machine learning models. Since this system is implemented by comparative analysis of different algorithms such as random forest algorithm, Extra tree classifier algorithm and XG Booster algorithm, it ensures high accuracy of predictions, in addition it allows the system to capture complex relationships and patterns in the data. The system aids in designing pollution control measures like issuing timely alerts and optimizing resource allocation for air quality management in India. It utilizes machine learning algorithms for the analysis and prediction of air pollution. The system's ability to capture complex relationships and provide accurate forecasts can contribute to mitigate the adverse effects of air pollution and improve public health.

### V Methodology

A deployed full-stack web application is implemented alongside integrating HTML, CSS, Bootstrap, JavaScript and Django with machine learning algorithms such as KNN, Random Forest Classifier, Decision Tree Classifier and Extra Trees Classifier. This methodology encompasses data collection, preprocessing, feature selection, model training and evaluation.

### VI Module Description

**i) Data Preprocessing**

The Data Preprocessing module is pivotal for preparing air quality data for analysis. It starts with aggregating data from diverse sources followed by cleaning to remove inconsistencies and outliers. This stage also involves transforming the data into a standardized format ensuring compatibility with the machine learning models. Feature engineering is then applied to select relevant features, such as pollutant levels and weather conditions. The data is further processed for analysis resolution, this preparation is essential for accurate prediction of air pollution levels setting a strong foundation for analysis and model effectiveness. The data values which poses null values or Nan are

either dropped or it is replaced by the mean value of the data values in that specific column. The data values which are of object type are parsed to numeric values using label encoder that helps in using the mathematical calculations in the visualization and analysis phase.

## ii) Data Visualization

The Data Visualization module plays a pivotal role by transforming complex air quality datasets into easily interpretable visuals, such as graphs and heatmaps. These visualizations reveal trends and patterns in pollution across different regions and times, assisting in data analysis and model refinement. By making data more accessible, this module supports informed decision-making and public awareness, significantly contributing to the system's overall effectiveness in predicting air pollution levels

## iii) Extra tree classifier Algorithm Implementation

The Extra Tree Classifier algorithm implementation module in this study is a critical component designed to enhance the predictive accuracy of air pollution levels in India. Leveraging an ensemble of extremely randomized trees, this module differentiates itself by the random selection of cut-points and features at each split in the decision tree construction process, thereby reducing the variance of the model without significantly increasing bias. This approach allows for a more robust model by integrating multiple decision trees to mitigate the risk of overfitting to the training data. In the context of air pollution prediction, the Extra Tree Classifier effectively captures complex nonlinear relationships between various predictors such as meteorological conditions, pollutant concentrations and temporal variables. The implementation involves feeding preprocessed and normalized data into the algorithm, followed by tuning hyperparameters to optimize performance. The outcome is a predictive model that offers significant insights into pollution trends and hotspot identification, making it an invaluable tool for environmental monitoring and decision-making in mitigating air pollution in India.

## iv) Implementing Random Forest Tree Algorithm

In the implementation of the Random Forest algorithm module, this approach focuses on leveraging this ensemble learning method for its robustness and accuracy in predicting air pollution levels. Random Forest operates by constructing multiple decision trees during the training phase and outputting the mode of the classes (classification) or mean prediction (regression) of the individual trees. This method is particularly suitable for the dataset because it can handle the high variability and complexity of environmental data, including pollutant concentrations, meteorological factors and temporal patterns. By integrating the Random Forest algorithm, we aim to improve the predictive accuracy of the system, reduce overfitting risks inherent in decision trees and provide a reliable tool for forecasting air pollution. This module processes preprocessed data, applies feature selection techniques to identify the most influential variables and then trains the Random Forest model using a split of training and testing datasets to evaluate its performance. The outcome is a robust model that enhances the system's ability to predict air pollution levels with high accuracy, making it a vital component of overall architecture aimed at tackling air quality issues in India.

## v) Implementing XG Boost Classifier Algorithm

The XG Boost Classifier algorithm implementation module in this paper focuses on integrating the XGBoost (Extreme Gradient Boosting) algorithm into the air pollution prediction system. XGBoost is a powerful machine learning algorithm known for its efficiency and accuracy in handling large datasets and complex relationships. This module involves preprocessing the data to ensure compatibility with XGBoost, tuning hyperparameters to optimize model performance, training the classifier on historical air quality data and evaluating its effectiveness using appropriate metrics such as accuracy, precision and recall. Additionally, the module includes techniques for feature importance analysis to gain insights into the factors driving air pollution levels. By incorporating XGBoost into the system, we aim to enhance the predictive capabilities and robustness of the models ultimately improving the accuracy of air pollution forecasts and supporting informed decision-making for environmental management.

## vi) Algorithm Analysis and Selection

Selecting the best algorithm depends on the following performance metrices:

**Accuracy:** Accuracy measures the proportion of correctly classified instances out of the total instances. It is calculated as the ratio of the number of correct predictions to the total number of predictions. While accuracy is a widely used metric, it may not be suitable for imbalanced datasets where one class dominates the others.

**Precision:** Precision measures the proportion of true positive predictions out of all positive predictions. It focuses on the accuracy of positive predictions and is calculated as the ratio of true positives to the sum of true positives and false positives. Precision is useful when the cost of false positives is high.
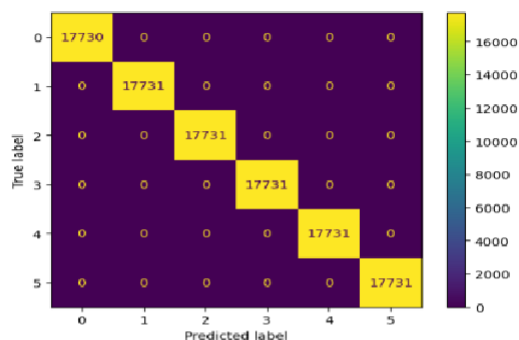
**Recall (Sensitivity):** Recall measures the proportion of true positive predictions out of all actual positive instances. It focuses on the ability of the model to capture all positive instances and is calculated as the ratio of true positives to the sum of true positives and false negatives. Recall is important when the cost of false negatives is high.

**F1 Score:** The F1 score is the harmonic mean of precision and recall, providing a balance between the two metrics. It is calculated as 2 * (precision * recall) / (precision + recall) and ranges from 0 to 1, with higher values indicating better performance.

**Specificity:** Specificity measures the proportion of true negative predictions out of all actual negative instances. It is the complement of the false positive rate and is useful in binary classification tasks with imbalanced classes.

The algorithms such as Extra tree classifier, random forest classifier and XG boost classifier are trained, implemented and the results are tested and analyzed to choose the one that composes high accuracy levels of predicting the exact output for a desired set of input from the testing data set.

The image below is the confusion matrix of the accuracy level in random forest classifier algorithm:

Considering the performance metrices, random forest classifier algorithm produces an accuracy level of 99.98812003571932.

### vii) Deployment using Django

The Deployment using Django implementation module facilitates the seamless integration and deployment of the developed air pollution analysis and prediction system. Leveraging the Django web framework, this module configures the web server environment, setting up the Django project structure, managing static and media files and ensures database connectivity. It automates the deployment process, streamlining tasks such as server configuration, database migration and application deployment. Security considerations are paramount, with measures implemented to safeguard data integrity and protect against potential threats. Continuous monitoring and logging mechanisms are established to track application performance and detect any issues. Additionally, user authentication and authorization mechanisms are implemented to control access to the deployed application, ensuring secure and controlled usage. Overall, the deployment using Django module ensures the efficient deployment and secure accessibility of the air pollution analysis and prediction system contributing to effective environmental management and public health outcomes.

## VI. Conclusion

The applications of machine learning for predicting air pollution holds immense promise to address the critical issue of air quality in urban environment. The machine learning models have shown the ability to accurately predict pollution levels, identify pollution sources and offer necessary actions or recommendations for mitigation.

## VII Future Work

Advanced Model Interpretability: It will focus more on developing advanced techniques for interpreting machine learning models that are used to predict air pollution. It helps us to understand various factors for air pollution and improve the transparency of model recommendations.

Integration with Environmental Policies: It will explore variouss ways to seamlessly integrate machine learning predictions into environmental policies and regulations thereby ensuring that data-driven insights directly inform decision-making and pollution reduction strategies.

## VIII. References

[1] Sai Bhargav Kasetty, S. Nagini A Survey Paper on an IoT-based Machine Learning Model to Predict Air Pollution Levels 2022 4th International Conference on Advances in Computing, Communication Control and Networking (ICAC3N) Year: 2022 | Conference Paper | Publisher: IEEE.

[2] Peijiang Zhao, Koji Zettsu Convolution Recurrent Neural Networks Based Dynamic Transboundary Air Pollution Predictiona 2019 IEEE 4th International Conference on Big Data Analytics (ICBDA) Year: 2019 | Conference Paper | Publisher: IEEE

[3] Shahan Salim, Irfhana Zakir Hussain, Jasleen Kaur, Plinio P. Morita "An Early Warning System for Air Pollution Surveillance: A Big Data Framework to Monitoring Risks Associated with Air Pollution" 2023 IEEE International Conference on Big Data (BigData) Year: 2023 | Conference Paper | Publisher: IEEE

[4] Nurul Aini, M Syukri Mustafa "Data Mining Approach to Predict Air Pollution in Makassar" 2020 2nd International Conference on Cybernetics and Intelligent System (ICORIS) Year: 2020 | Conference Paper | Publisher: IEEE

[5] Amar Catovic, Esad Kadusic, Christoph Ruland, Natasa Zivic, Narcisa Hadzajlic "Air pollution prediction and warning system using IoT and machine learning" 2022 International Conference on Electrical, Computer, Communications and Mechatronics Engineering (ICECCME) Year: 2022 | Conference Paper | Publisher: IEEE

[6] Ankeshit Srivastava, Ayaz Ahmad, Sunny Kumar, Md Arman Ahmad "Air Pollution Data and Forecasting Data Monitored through Google Cloud Services by using Artificial Intelligence and Machine Learning" 2022 6th International Conference on Electronics, Communication and Aerospace Technology Year: 2022 | Conference Paper | Publisher: IEEE

[7] Harshit Srivastava, Goutam Kumar Sahoo, Santos Kumar Das, Poonam Singh "Performance Analysis of Machine Learning Models for Air Pollution Prediction" 2022 International Conference on Smart Generation Computing, Communication and Networking (SMART GENCON) Year: 2022 | Conference Paper | Publisher: IEEE

[8] Harshal P. Varade, Sonal C. Bhangale, Sandip R. Thorat, Pravin B. Khatkale, Santosh Kumar Sharma, P. William "Framework of Air Pollution Assessment in Smart Cities using IoT with Machine Learning Approach" 2023 2nd International Conference on Applied Artificial Intelligence and Computing (ICAAIC) Year: 2023 | Conference Paper | Publisher: IEEE

[9] Madhushika Mihirani, Lasith Yasakethu, Sachintha Balasooriya "Machine Learning-based Air Pollution Prediction Model" 2023 IEEE IAS Global Conference on Emerging Technologies (GlobConET) Year: 2023 | Conference Paper | Publisher: IEEE

[10] Usha Mahalingam, Kirthiga Elangovan, Himanshu Dobhal, Chocko Valliappa, Sindhu Shrestha, Giriprasad Kedam "A Machine Learning Model for Air Quality Prediction for Smart Cities" 2019 International Conference on Wireless Communications Signal Processing and Networking (WiSPNET) Year: 2019 | Conference Paper | Publisher: IEEE

# RE-2022-221164-plag-report

| **17**% | **15**% | **11**% | **12**% |
|---|---|---|---|
| SIMILARITY INDEX | INTERNET SOURCES | PUBLICATIONS | STUDENT PAPERS |

PRIMARY SOURCES

**1**    **Submitted to Liverpool John Moores University**    **2**%
Student Paper

**2**    **Submitted to Heriot-Watt University**    **1**%
Student Paper

**3**    **www.ijritcc.org**    **1**%
Internet Source

**4**    **tylervigen.com**    **1**%
Internet Source

**5**    **Submitted to Manipal University**    **1**%
Student Paper

**6**    **Submitted to Swinburne University of Technology**    **1**%
Student Paper

**7**    **Submitted to University of Auckland**    **1**%
Student Paper

**8**    **Submitted to Vilnius Gediminas Technical University**    **1**%
Student Paper

| 9 | www.researchgate.net<br>Internet Source | 1% |
|---|---|---|
| 10 | Submitted to CSU, San Jose State University<br>Student Paper | 1% |
| 11 | files.eric.ed.gov<br>Internet Source | <1% |
| 12 | ijritcc.org<br>Internet Source | <1% |
| 13 | Submitted to Monash University<br>Student Paper | <1% |
| 14 | M. Bindhu, S. Parasuraman, S. Yogeeswaran, S. Nimmi Devi. "Harnessing Machine Learning for IoT-Driven Atmospheric Parameter Monitoring and Predictive Analytics", 2023 9th International Conference on Smart Structures and Systems (ICSSS), 2023<br>Publication | <1% |
| 15 | Submitted to Silpakorn University<br>Student Paper | <1% |
| 16 | arxiv.org<br>Internet Source | <1% |
| 17 | hess.copernicus.org<br>Internet Source | <1% |
| 18 | Antonio Pagliaro. "Forecasting Significant Stock Market Price Changes Using Machine | <1% |

Learning: Extra Trees Classifier Leads",
Electronics, 2023
Publication

| 19 | Submitted to University of Essex<br>Student Paper | <1% |
|---|---|---|

| 20 | eprints.utar.edu.my<br>Internet Source | <1% |
|---|---|---|

| 21 | i-scholar.in<br>Internet Source | <1% |
|---|---|---|

| 22 | Submitted to Harrisburg University of Science and Technology<br>Student Paper | <1% |
|---|---|---|

| 23 | Kurt, A.. "Forecasting air pollutant indicator levels with geographic models 3days in advance using neural networks", Expert Systems With Applications, 201012<br>Publication | <1% |
|---|---|---|

| 24 | N S Aruna Kumari, K S Ananda Kumar, S Hitesh Vardhan Raju, H R Vasuki, M P Nikesh. "Prediction of Air Quality in Industrial Area", 2020 International Conference on Recent Trends on Electronics, Information, Communication & Technology (RTEICT), 2020<br>Publication | <1% |
|---|---|---|

| 25 | Sudhakar Pal, Arabinda Sharma. "How does the COVID-19-related restriction affect the spatiotemporal variability of ambient air | <1% |
|---|---|---|

quality in a tropical city?", Environmental Monitoring and Assessment, 2023
Publication

| 26 | worldwidescience.org | <1% |
| --- | --- | --- |
| | Internet Source | |

| 27 | www.frontiersin.org | <1% |
| --- | --- | --- |
| | Internet Source | |

| 28 | www.mdpi.com | <1% |
| --- | --- | --- |
| | Internet Source | |

| 29 | www.science.gov | <1% |
| --- | --- | --- |
| | Internet Source | |

| 30 | Shomya Kumari, Deepak Kumar, Manish Kumar, Chaitanya B. Pande. "Modeling of standardized groundwater index of Bihar using machine learning techniques", Physics and Chemistry of the Earth, Parts A/B/C, 2023 | <1% |
| --- | --- | --- |
| | Publication | |

Exclude quotes          On          Exclude matches          Off
Exclude bibliography   On