# Air Pollution Forecasting in NCR

SUBMITTED IN PARTIAL FULFILLMENT FOR THE REQUIREMENT OF THE
AWARD OF DEGREE OF

## BACHELOR OF TECHNOLOGY

## IN

## COMPUTER SCIENCE



Submitted by

Avi Chaudhary (2000290120048)

Avika Tyagi (2000290120049)

Ankita Kushwaha (2000290120026)

.

Supervised by:

Prof. Abhishek Goyal

**Session – 2023-24**
DEPARTMENT OF COMPUTER SCIENCE
## KIET GROUP OF INSTITUTIONS, GHAZIABAD

**(Affiliated to Dr. A. P. J. Abdul Kalam Technical University, Lucknow, U.P., India)**

# DECLARATION

We hereby declare that this submission is our own work and that, to the best of our knowledge and belief, it contains no material previously published or written by another person nor material which to a substantial extent has been accepted for the award of any other degree or diploma of the university or other institute of higher learning, except where due acknowledgment has been made in the text.

Signature

Name:-

Roll No.:-

Date:- 22/04/2024

# CERTIFICATE

This is to certify that Project Report entitled "Air Pollution Forecasting in Delhi NCR" which is submitted by Avi Chaudhary,Avika Tyagi, and Ankita Kushwaha in partial fulfillment of the requirement for the award of degree B. Tech. in Department of Computer Science of Dr. A.P.J. Abdul Kalam Technical University, Lucknow is a record of the candidates own work carried out by them under my supervision. The matter embodied in this report is original and has not been submitted for the award of any other degree.

.

**Date: 22/04/2024**                                          **Supervisor**

Name – Prof.

Abhishek Goyal

ACKNOWLEDGEMENT

It gives us a great sense of pleasure to present the report of the B. Tech Project undertaken during B. Tech. Final Year. We owe special debt of gratitude to Professor Abhishek Goyal, Department of Computer Science, KIET, Ghaziabad, for his constant support and guidance throughout the course of our work. Her sincerity, thoroughness and perseverance have been a constant source of inspiration for us. It is only his cognizant efforts that our endeavors have seen light of the day.

We also take the opportunity to acknowledge the contribution of Dr. Ajay Kumar Srivastava, Head of the Department of Computer Science, KIET, Ghaziabad, for his full support and assistance during the development of the project. We also do not like to miss the opportunity to acknowledge the contribution of all the faculty members of the department for their kind assistance and cooperation during the development of our project.

Last but not the least, we acknowledge our friends for their contribution in the completion of the project.

Date:

Signature:                                          Signature:

Name: Avi Chaudhary                      Name: Avika Tyagi

Roll No.: 2000290120048                Roll No.: 2000290120049

Signature:

Name: Ankita Kushwaha

Roll No.:2000290120026

# ABSTRACT

The escalation of urbanization and industrial activities in the National Capital Region (NCR) has led to growing concerns about air quality and pollution. In response, this project focuses on developing an advanced pollution forecasting system tailored specifically for the NCR region. Leveraging state-of-the-art technologies, including machine learning algorithms, the system analyzes historical air quality data, meteorological parameters, and other relevant factors to predict air quality levels accurately. By integrating data from diverse sources, the system addresses the complexity of pollution factors such as vehicular emissions, industrial processes, and meteorological influences. Through a comprehensive comparative study evaluating various metrics like accuracy and precision under different scenarios, the project lays the foundation for continuous improvement in pollution forecasting models. The system's graphical presentation of results showcases its performance across different parameters, including prediction accuracy, response time, and adaptability to changing environmental conditions. Positioned as a critical tool for environmental management in the NCR, the pollution forecasting system aids decision-making related to public health and regulatory interventions. Insights gleaned from the project guide future research efforts to enhance the system's accuracy and scope, with ongoing investigations including factors like traffic congestion and dynamic emission sources for more comprehensive pollution modeling.

# TABLE OF CONTENTS

# CHAPTER 1
# INTRODUCTION

## 1.1   INTRODUCTION

The National Capital Region (NCR) is grappling with a formidable adversary: the relentless surge in air pollution levels. This burgeoning crisis stems from a confluence of factors, including the rapid pace of urbanization, the proliferation of industrial activities, and the relentless march of vehicular emissions. As the skyline becomes increasingly obscured by smog and particulate matter, the urgency to combat this environmental menace becomes ever more palpable.

In response to this looming threat, this project emerges as a beacon of hope and innovation. It embarks on the mission to craft an advanced pollution forecasting system meticulously tailored to the intricate dynamics of the NCR. Beyond merely acknowledging the gravity of the situation, the project sets out to forge a solution that transcends conventional boundaries, drawing upon cutting-edge technologies and interdisciplinary approaches.

At its core, the overarching goal of this endeavor is crystal clear: to furnish the NCR with a potent arsenal in its battle against air pollution. By harnessing the power of predictive analytics, machine learning algorithms, and real-time data streams, the envisioned system aspires to be more than a passive observer; it aims to be a proactive guardian of the environment, equipped with the foresight to anticipate pollution spikes and mitigate their impact.

In essence, this project represents a pivotal step forward in the realm of environmental stewardship. It underscores the imperative to embrace innovation as a catalyst for change and underscores the imperative to embrace innovation as a catalyst for change and resilience in the face of environmental adversity. By empowering stakeholders with actionable insights and foresight, it seeks to catalyze a paradigm shift towards a future where clean air is not merely an aspiration but a tangible reality for all residents of the NCR."

## 1.2   PROJECT CATEGORY

The provided content falls under the category of "Air Pollution Forecasting and Management Using Machine Learning." This project is centered around developing an advanced pollution forecasting system tailored specifically for the National Capital Region (NCR). By integrating state-of-the-art technologies and employing machine learning algorithms like Support Vector Machine (SVM), K-Nearest Neighbors (KNN), and Multiple Linear Regression (MLR), the project aims to predict air quality levels accurately. Visualization of data through graphs plays a crucial role in analyzing resource distribution and optimizing forecasting models. The project underscores the significance of achieving optimal conditions for air quality management in the NCR by leveraging predictive analysis and machine learning techniques. Through comprehensive research and data-driven approaches, it seeks to provide timely and effective solutions to combat air pollution and safeguard public health and environmental well-being in the region.

# CHAPTER 2
# LITERATURE REVIEW

## 2.1 LITERATURE REVIEW

**Anikender Kumar, Pramila Goyal** presented the study that forecasts the daily AQI value for the city        Delhi, India using previous record of AQI and meteorological parameters with the help of Principal    Component Regression (PCR) and Multiple Linear Regression Techniques. They perform the prediction of daily AQI of the year 2006 using previous records of the year 2000-2005 and different equations. After that this predicted value is then compared with the observed value of AQI of 2006 for the seasons summer, Monsoon, Post Monsoon and winter using Multiple Linear Regression Technique [1]. Principal Component Analysis is used to find the collinearity among the independent variables. The principal components were used in Multiple Linear Regression to eliminate collinearity among the predictor variables and also reduce the number of predictors [1]. The Principal Component Regression gives better performance for predicting the AQI in the winter season than any other seasons. In this study only meteorological parameters were considered or used while forecasting the future AQI but they have not considered the ambient air pollutants that may cause the adverse health effects.

**Huixiang Liu (et al.2019)** have taken two different cities Beijing and Italian city for the study purpose. They have forecasted the Air Quality Index (AQI) for the city Beijing and predicting the concentration of NOx in an Italian City depending on two different publicly available datasets. The first Dataset for the period of December 2013 to August 2018 having 1738 instances is made available from the Beijing Municipal Environmental Centre [5] which contains the fields like hourly averaged AQI and the concentrations of PM2.5, O3, SO2, PM10, and NO2 in Beijing. The second Dataset with 9358 instances is collected from Italian city for the period of March 2004 to February 2005. This dataset contains the attributes as Hourly averaged concentration of CO, non-methane Hydrocarbons, Benzene, NOx, NO2 [5]. But they focused majorly on NOx prediction as it is one of the important predictors for Air Quality evaluation.

**S. Ghude et al** Ghude and colleagues undertook a comprehensive investigation into the realm of air quality forecasting, focusing their attention specifically on the bustling Indian metropolis of Mumbai. Recognizing the unique challenges posed by the city's densely populated urban landscape and its intricate interplay of meteorological factors and pollution sources, the researchers embarked on a quest to develop an innovative forecasting framework tailored to Mumbai's distinctive characteristics. Drawing upon a rich repository of historical air quality data spanning several years, they employed sophisticated machine learning techniques, including Artificial Neural Networks (ANNs) and Decision Trees, to unravel the complex patterns underlying air pollution dynamics. Central to their methodology was the integration of local meteorological variables, such as temperature, humidity, wind speed, and atmospheric stability, into the forecasting models, thereby capturing the nuanced interactions between atmospheric conditions and pollutant dispersion. By leveraging the power of data-driven modeling approaches, Ghude and his team sought to enhance the accuracy and reliability of air quality predictions, laying the groundwork for proactive pollution mitigation strategies and informed decision-making in Mumbai's quest for cleaner, healthier air.

**Y. Li et al. (2018**)   Li and his research team dedicated their efforts to refining air quality forecasting methodologies tailored for urban landscapes. With rapid urbanization and industrialization exacerbating air pollution levels in cities worldwide, their study sought to enhance forecasting accuracy through an innovative hybrid approach. By integrating data-driven models with physical models, their research aimed to leverage the strengths of both approaches while mitigating their respective limitations. Drawing upon a diverse array of data sources, including satellite observations, meteorological data, and ground-based monitoring data, their methodology aimed to capture the intricate interplay between various atmospheric parameters and pollutant concentrations. Through meticulous validation and testing, their study demonstrated the efficacy of the hybrid approach in generating precise and reliable air quality forecasts, thereby empowering policymakers and urban planners with invaluable insights for sustainable development and environmental management.

**Zhang et al. (2017).** Zhang and co-authors embarked on a nuanced investigation into air quality forecasting, spotlighting the nexus between transportation emissions and urban air pollution. As urbanization accelerates and vehicular traffic burgeons in cities worldwide, their study sought to elucidate the complex interactions between transportation activities and air quality dynamics. By

analyzing comprehensive datasets encompassing vehicular traffic patterns, road infrastructure, emission inventories, and air quality measurements, their research aimed to uncover the underlying mechanisms driving pollution levels in urban environments. Through sophisticated modeling techniques and statistical analyses, their study revealed the significant impact of transportation-related factors on air quality, underscoring the need for holistic strategies to mitigate emissions and improve urban air quality. The findings from their research not only provided valuable insights into the factors shaping air pollution in urban areas but also informed the development of targeted interventions and policies aimed at fostering sustainable urban mobility and reducing environmental burdens.

**X. Wang et al. (2020)** Wang and collaborators delved into the realm of air quality forecasting armed with cutting-edge sensor technologies tailored for urban environments. Recognizing the limitations of traditional monitoring approaches in capturing real-time air quality dynamics, their research sought to harness the potential of advanced sensor networks for enhanced forecasting capabilities. By deploying a network of low-cost sensors capable of real-time pollutant detection, their study aimed to overcome spatial and temporal limitations inherent in conventional monitoring systems. Through a combination of sensor data integration and machine learning algorithms, their methodology aimed to generate high-resolution air quality forecasts at a local scale, providing stakeholders with timely and actionable information for pollution mitigation efforts. The findings from their research not only demonstrated the feasibility of sensor-based forecasting approaches but also highlighted the transformative potential of emerging technologies in revolutionizing air quality monitoring and management in urban environments.

**C. Wei et al. (2019)** Wei and his research team embarked on a meticulous examination of the intricate interplay between meteorological conditions and air quality variations in a coastal city. Recognizing the profound influence of meteorological factors on pollutant dispersion and atmospheric dynamics, their study sought to elucidate the complex relationships shaping air quality patterns in coastal urban environments. By conducting a comprehensive analysis of atmospheric dynamics, including temperature inversions, sea breeze effects, and atmospheric stability, their research aimed to unravel the underlying mechanisms driving air pollution episodes. Through sophisticated modeling techniques and statistical analyses, their study revealed the multifaceted interactions between meteorology and air quality, underscoring the need for tailored forecasting

approaches that account for localized atmospheric dynamics. The findings from their research not only enhanced our understanding of the factors shaping air quality in coastal cities but also provided valuable insights for developing effective pollution mitigation strategies and resilience measures in the face of changing climate conditions.

**M. Kumar et al. (2018).** Kumar and colleagues delved into the role of urban green spaces as a potential panacea for mitigating air pollution and enhancing air quality. Recognizing the mounting concerns regarding deteriorating air quality in urban areas worldwide, their study sought to explore nature-based solutions for pollution mitigation. By conducting a combination of field studies and modeling experiments, their research aimed to quantify the impact of vegetation cover on pollutant removal and dispersion dynamics. Through meticulous analysis of vegetation characteristics, pollutant deposition rates, and atmospheric dispersion patterns, their study revealed the critical role of urban green spaces in enhancing air quality and promoting public health. The findings from their research not only underscored the multifaceted benefits of green infrastructure but also advocated for its integration into urban planning and environmental policymaking frameworks as a sustainable solution for addressing air pollution challenges in urban environments.

**H. Chen et al. (2019)** Chen and his research team embarked on a quest to assess the utility of air quality forecasting systems in informing public health interventions and policy decisions. Recognizing the imperative of proactive measures to mitigate the adverse health effects of air pollution, their study sought to evaluate the effectiveness of forecasting models in predicting pollutant concentrations and associated health risks. By conducting a comparative analysis of different forecasting techniques and validation against observational data, their research aimed to elucidate the strengths and limitations of existing forecasting systems. Through rigorous evaluation metrics and sensitivity analyses, their study provided valuable insights into the performance of forecasting models under varying environmental conditions and pollutant sources. The findings from their research not only underscored the importance of integrating air quality forecasts into decision-making processes but also informed the development of targeted interventions and policy measures aimed at safeguarding public health and promoting environmental sustainability.

## 2.2 PROBLEM FORMULATION

The project's foundation lies in addressing the multifaceted challenges posed by the surge in air pollution within the NCR. The formulation of the problem revolves around the intricate interactions of diverse pollution contributors, including vehicular emissions, industrial processes, and the influence of meteorological conditions. Understanding these complex interactions is crucial for devising effective strategies to mitigate air pollution and safeguard public health. By delving into the intricate web of factors influencing air quality, the project aims to develop a comprehensive understanding of the underlying mechanisms driving pollution levels in the NCR. This holistic approach enables the project to account for the dynamic nature of air pollution and its spatial and temporal variability across the region. Moreover, by integrating advanced data analysis techniques and modeling methodologies, the project seeks to unravel the causal relationships between different pollution sources and their cumulative impact on air quality. Through interdisciplinary collaboration and stakeholder engagement, the project endeavors to forge innovative solutions that transcend traditional boundaries.

## 2.3  OBJECTIVES

**Develop an Advanced Pollution Forecasting System:**

Craft a cutting-edge system capable of accurately predicting pollution levels in the NCR through the utilization of advanced modeling techniques and state-of-the-art technology.

**Integration of Data Sources:**

Integrate diverse data sources, including historical air quality data, meteorological parameters, and real-time information, to enhance the accuracy and reliability of pollution forecasts by capturing the full spectrum of influencing factors.

**Utilization of Machine Learning Algorithms:**

Leverage sophisticated machine learning algorithms to model the complex interactions between various pollution contributors, enabling a holistic and nuanced approach to pollution forecasting that accounts for dynamic environmental dynamics.

**Conduct Comparative Studies:**

Undertake comprehensive comparative studies to assess the performance of the forecasting system under varied scenarios, providing invaluable insights into its efficacy and areas for enhancement.

**Empower Stakeholders:**

Equip stakeholders with actionable insights derived from the forecasting system, enabling informed decision-making and proactive environmental management strategies to mitigate the adverse effects of air pollution.

**Design for Scalability:**

Develop a scalable model capable of adapting to evolving pollution patterns in the NCR, ensuring the long-term effectiveness and sustainability of the forecasting system as environmental dynamics evolve.

**Contribute to Scientific Understanding:**

Advance scientific understanding of pollution dynamics in urban environments by contributing valuable insights and data-driven analyses, fostering ongoing research and innovation in the field of environmental science and management.

# CHAPTER 3
# PROPOSED SYSTEM

## 3.1 PROPOSED SYSTEM

The proposed system is not just a standalone solution but a dynamic and adaptable framework that evolves with the changing environmental landscape of the NCR. By integrating data from diverse sources, including ground-based sensors, satellite observations, and meteorological stations, it ensures a holistic understanding of pollution dynamics. Furthermore, the utilization of machine learning algorithms empowers the system to continuously learn and improve its predictive capabilities over time, enhancing its effectiveness in addressing the complex challenges posed by air pollution. Through ongoing refinement and validation, the system strives to become a cornerstone in environmental management efforts, facilitating evidence-based decision-making and proactive measures to safeguard public health and the environment

## 3.2 UNIQUE FEATURES OF THE SYSTEM

**Tailored Forecasting:**

One of the most remarkable aspects of the system lies in its capacity to offer tailored predictions, catering to the diverse array of pollution sources prevalent in the NCR. Rather than employing a one-size-fits-all approach, the system acknowledges the unique characteristics and contributions of each pollutant source, allowing for nuanced and accurate forecasting tailored to specific environmental conditions and emission profiles.

**Integrated Data Architecture:**

At the heart of the system's functionality is its seamless integration of data from various sources, including historical records, meteorological parameters, and real-time observations. This integrated approach not only enhances the system's predictive accuracy by providing a comprehensive understanding of pollution dynamics but also fosters synergy between different data streams, facilitating holistic analysis and informed decision-making.

**Machine Learning Algorithms:**

Central to the system's analytical prowess are its sophisticated machine learning algorithms, which enable it to decipher complex interactions and patterns within the vast expanse of environmental data. By harnessing the power of these advanced algorithms, the system can uncover hidden correlations, adapt to changing conditions, and continuously refine its predictive models, ensuring that forecasts remain accurate and reliable even in the face of evolving environmental challenges.

**Scalability:**

An essential characteristic of the system is its scalability, allowing it to grow and adapt alongside the dynamic nature of pollution patterns and environmental conditions in the NCR. By designing the system with scalability in mind, stakeholders can be confident in its ability to accommodate increasing data volumes, emerging technologies, and evolving pollution mitigation strategies, thereby ensuring its long-term relevance and effectiveness in addressing air quality challenges.

# CHAPTER 4

# REQUIREMENT ANALYSIS AND SYSTEM SPECIFICATION

## 4.1   FEASIBILITY STUDY

Conducting a feasibility study for an automated retail checkout system using Algorithms involves evaluating its technical, economic, and operational aspects. Following are the various components:

### 4.1.1 TECHNICAL FEASIBILITY:

Technical feasibility evaluates the project's ability to be implemented from a technological standpoint. This involves assessing existing infrastructure, compatibility with required technology, and the availability of necessary resources. For our project, we need to ensure that the machine learning algorithms and data processing techniques can be effectively integrated into the 5G network environment. Compatibility with existing network components and data sources will also be analyzed.

#### 4.1.1.1 System Requirements:

the system must be capable of processing and analyzing data in real-time, considering factors like application type, signal strength, and latency to optimize bandwidth allocation. The availability of necessary resources, including computational power and data storage, will also be crucial for the successful implementation of the project. Therefore, technical feasibility will be assessed by evaluating the project's compatibility with existing technology, infrastructure, and resource availability within the 5G network environment.

### 4.1.2 ECONOMIC FEASIBILITY:

#### 4.1.2.1 Cost Analysis:

The cost analysis for air pollution forecasting in the National Capital Region (NCR) involves a comprehensive assessment of both initial and ongoing expenses associated with the project

implementation. Initial costs encompass hardware and software acquisition, development and integration expenses, training costs, and any other upfront investments necessary for deploying machine learning models and infrastructure. Ongoing costs include maintenance, monitoring, and operational expenses incurred over the project's lifecycle. By meticulously evaluating these costs, stakeholders can gain insights into financial commitments, enabling informed decisions on resource allocation and budgeting to ensure cost-effectiveness and sustainability.

### 4.1.2.2. Return on Investment (ROI):

ROI analysis evaluates the financial benefits relative to the investment cost. For this project, ROI involves calculating the monetary returns from implementing machine learning models (such as SVM and KNN) for air pollution forecasting in the NCR. It includes estimating the initial investment required for technology integration and data processing tools against expected benefits like improved air quality monitoring and potential cost savings from timely interventions. By comparing projected returns with investment costs, ROI analysis provides insights into project profitability and efficiency.

### 4.1.2.3. Payback Period:

The payback period signifies the time taken for an investment to recover its initial costs through generated returns. In this project's context, the payback period is determined by analyzing the timeline for realizing benefits from air pollution forecasting. Factors considered include adoption rate, effectiveness in pollution mitigation, and potential cost savings from preventive measures. Identifying the point where cumulative returns match or exceed the initial investment, the payback period aids in assessing financial feasibility and aligning with organizational goals for achieving favorable returns within a reasonable timeframe.

## 4.1.3 OPERATIONAL FEASIBILITY

**Availability of Technology and Expertise:** The success of air pollution forecasting in the National Capital Region (NCR) hinges on access to cutting-edge technologies like machine learning algorithms (SVM, KNN, MLR) and data visualization tools. Expertise in network engineering, data analytics, and machine learning is essential for integrating these technologies effectively. Access to skilled professionals and fostering collaboration between technical teams is crucial for leveraging these technologies to tackle the challenges posed by air pollution effectively.

**Robustness of Forecasting Models:** Machine learning algorithms such as SVM and KNN demonstrate robustness in predicting air pollution levels and optimizing forecasting accuracy. SVM excels in classification and regression tasks, while KNN offers instance-based learning for precise predictions. Their combined strengths ensure effective pollution forecasting and application consistency across diverse environmental conditions within the NCR.

**Scalability and Efficiency:** The design of air pollution forecasting systems must prioritize scalability and efficiency to adapt to varying pollution patterns and meteorological conditions. Leveraging machine learning models like SVM and KNN enhances scalability by efficiently predicting pollution trends and optimizing resource allocation, thereby maximizing overall system efficiency.

**Integration with Existing Systems:** Integrating machine learning models (SVM, KNN) into air pollution forecasting systems aims to optimize prediction accuracy and resource allocation seamlessly. Utilizing data visualization tools facilitates aligning forecasting results with existing environmental monitoring systems, enhancing the overall effectiveness of pollution mitigation strategies in the NCR.

**Flexibility and Adaptability:** The system exhibits flexibility and adaptability by dynamically adjusting forecasting parameters based on real-time data inputs and environmental variables. By

considering factors such as pollutant sources, meteorological conditions, and historical trends, the system ensures accurate and tailored predictions for different pollution contributors within the NCR, enabling proactive environmental management and public health interventions.

# 4.2 SOFTWARE REQUIREMENT SPECIFICATION DOCUMENT

## 4.2.1  DATA REQUIREMENT

**Data Sources**

In order to train SVM, KNN on custom dataset, we have first imported the model from the respective official repository. And further tuned their latent codes and weights to train the data accordingly.

**Data Set**

The dataset is analyzed to understand the impact of parameters like signal strength and latency on application performance. Visualized through graphs and charts, it guides resource distribution analysis, aiding in optimizing application performance.

## 4.2.2 FUNCTIONAL REQUIREMENT

**Data Requirements:**

The system must optimize latency for different types of applications. This encompasses a detailed specification of the types of data needed, including historical air quality data, meteorological parameters, and relevant contextual information.

**Machine Learning Model Integration:** The system should integrate three machine learning models—Support Vector Machine (SVM), K-Nearest Neighbors (KNN), and Multiple Linear Regression (MLR)—to predict and optimize bandwidth allocation based on various factors including application type, signal strength, and latency.

**Visual Representation of Data:** The system must provide visual representations of data through graphs and charts to facilitate resource distribution analysis, enabling users to understand and interpret network conditions and application performance based on signal strength and other metrics.

**Classification and Regression Capabilities:** SVM should excel in classification and regression tasks, while KNN should adopt instance-based learning methods to provide accurate predictions for bandwidth allocation. MLR should extend predictions by considering multiple variables to enhance the accuracy of predictions.

**Application Consistency and Performance Enhancement:** The system must ensure application consistency and performance enhancement by effectively allocating bandwidth based on predictive analysis and understanding of resource distribution, guided by SVM and KNN machine learning models.

## 4.2.2 PERFORMANCE REQUIREMENT

Performance requirements specify the performance metrics that the system must meet to ensure efficient operation. This includes metrics such as throughput, response time, and scalability. The system must be able to handle varying levels of network traffic and adapt dynamically to changing conditions while maintaining optimal performance.

## 4.2.3 MAINTAINABILITY REQUIREMENT

Maintainability requirements focus on ensuring that the system can be easily maintained and updated over time. This includes documentation, modular design, and support for future updates and enhancements. We need to design the system in a way that facilitates ongoing maintenance and allows for seamless integration of new features and improvements.

**Scalability:** The system should be designed to accommodate potential growth in the number of users. This scalability ensures that the infrastructure and machine learning models can handle increased demands without significant performance degradation.

**Modularity:** The system architecture should be modular, allowing for easy integration of new machine learning algorithms or updates to existing ones. Modularity facilitates flexibility and simplifies maintenance tasks by isolating components, making it easier to identify and address issues without disrupting the entire system.

**Documentation:** Comprehensive documentation should be provided for all aspects of the system, including the implementation of machine learning models, data processing techniques, and network infrastructure. This documentation aids in knowledge transfer, training new personnel, and troubleshooting, ensuring continuity in system maintenance and support.

**Version Control:** Implementing version control systems for code repositories and configuration files ensures that changes made to the system are tracked, logged, and reversible if necessary. Version control facilitates collaboration among team members, reduces the risk of errors during updates, and enables rollback to previous working states in case of issues.

**Error Logging and Monitoring:** The system should incorporate robust error logging and monitoring mechanisms to track system performance, identify anomalies, and proactively address potential issues. Real-time monitoring of network signals, application latency, and bandwidth utilization enables timely intervention to maintain optimal performance and user experience.

**Regular Updates and Maintenance:** Regular updates and maintenance activities should be scheduled to address software patches, security vulnerabilities, and performance optimizations.

## 4.2.5 SECURITY REQUIREMENT:

**Security Requirement:**

Security requirements address the safeguarding of sensitive environmental data and the integrity of the forecasting system This involves implementing measures to protect against unauthorized access, data breaches, and system vulnerabilities.

# 4.3 SDLC MODEL TO BE USED

The test methodology selected for the project is Agile. An Agile methodology is the most suitable for this project. It allows for flexibility, ongoing testing, and adaptation, which are essential for projects that involve machine learning and its algorithm. Agile enables you to respond to changing requirements and refine the style transfer algorithm as you gain insights from testing and user feedback. It is a sequential development process that flows like a waterfall through all phases of a project (analysis, design, development, and testing, for example), with each phase completely wrapping up before the next phase begins.

Following is the overview of all the phases of software development life cycle model:

**Requirement Analysis:** To create an effective air pollution forecasting system, a comprehensive requirement analysis is essential. This involves identifying key factors such as geographical location, sources of pollution, meteorological data, and historical trends. Additionally, stakeholders' needs, including government agencies, public health organizations, and the general public, must be considered to ensure the system meets diverse user requirements.

## System Design:

- **High-level design:**

Define the architecture for integrating machine learning models. Specify components for processing data and generating insights for bandwidth optimization.

- **Detailed design:**

Design algorithms for SVM, KNN, and MLR to predict bandwidth needs based on application characteristics and network conditions. Define visualization techniques for representing resource distribution and signal strength.

- **Implementation:**
- Develop algorithms and data processing modules for SVM, KNN, and MLR integration.

- Implement visualization tools for graphical representation of resource distribution.

## Testing:

- **Unit Testing:**

  Test individual components/modules of SVM, KNN, and MLR algorithms to ensure they function as expected.

- **Integration Testing:**

  Verify the interoperability and integration of machine learning models within the 5G network environment.

- **System Testing:**

  Evaluate the overall system functionality, including bandwidth optimization and application performance enhancement.

**Deployment:** Deploy the system in a test environment to ensure compatibility and stability. Roll out the system gradually to users, monitoring performance and addressing any issues that arise.

**Validation and Verification:** Validate the effectiveness of machine learning models in optimizing bandwidth allocation and enhancing application performance. Verify that the system meets the defined requirements and objectives.

**Maintenance and Support:** Provide ongoing maintenance and support to address any issues and ensure the continued functionality of the system. Incorporate updates and improvements based on user feedback and changing requirements.
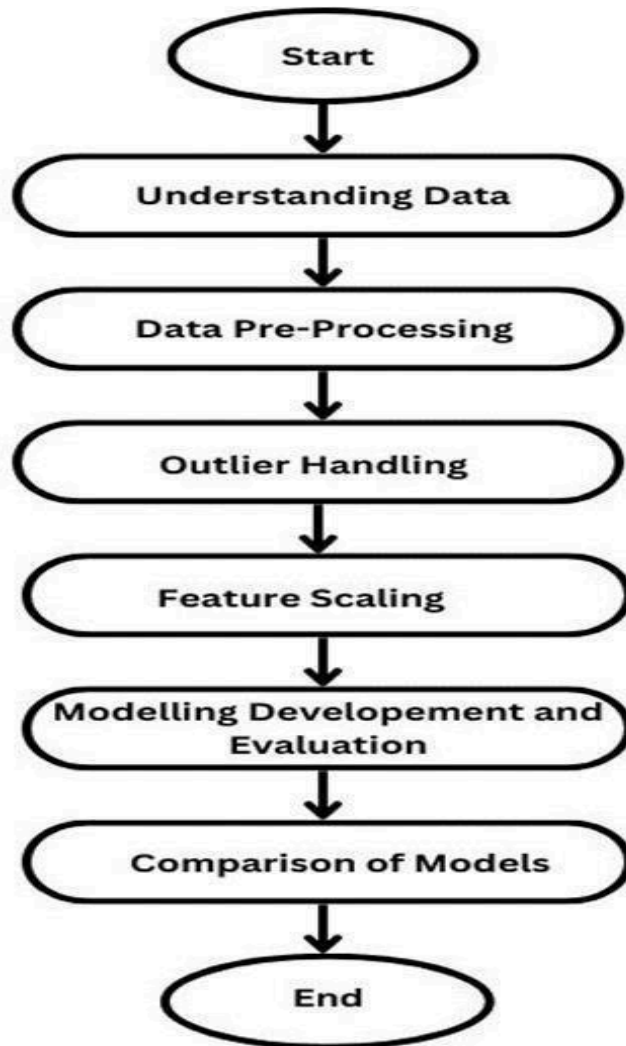
## 4.3.1 DETAIL DESIGN



**Fig. 3.1. Detailed Design of the System**

The data analysis pipeline typically begins with understanding and preprocessing data, including outlier handling and feature scaling. Then, modeling development and evaluation occur, ensuring compatibility and effectiveness.

## 4.3.2 SYSTEM DESIGN USING DFD LEVEL 0 AND LEVEL 1

**4.3.2.1** DFD Level 0



**Fig. 3.2. Level 0 DFD**

**4.3.2.2** DFD Level 1
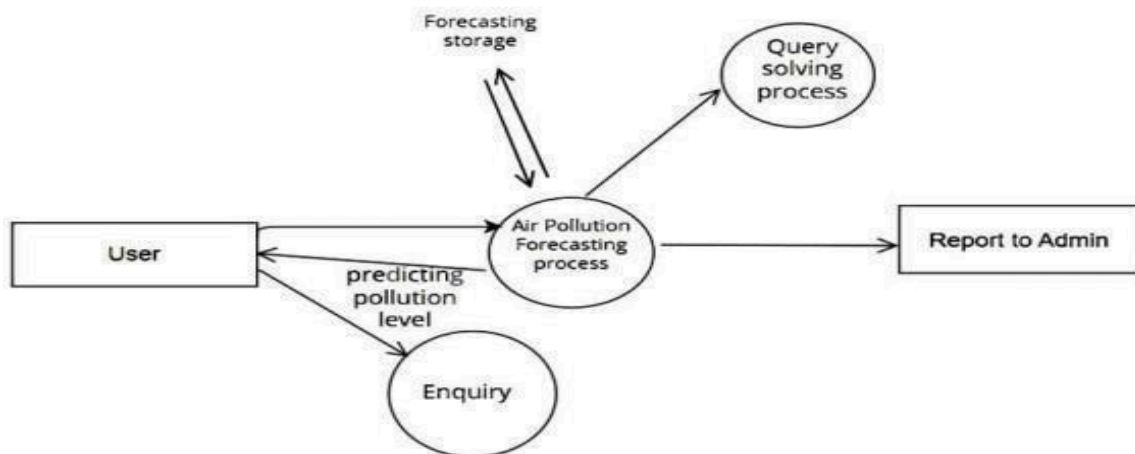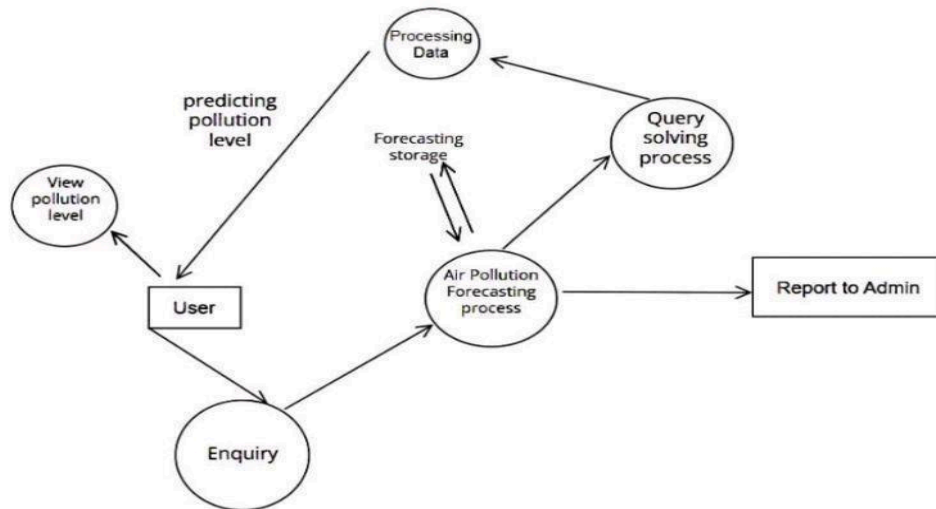


**Fig. 3.3. 1 Level 1 DFD**

## **4.3.2.3**   DFD Level 2



**Fig 3.4 Level 2 DFD**

## **4.3.2.4**  USE CASE DIAGRAM



**Fig. 3.8 Use Case Diagram**

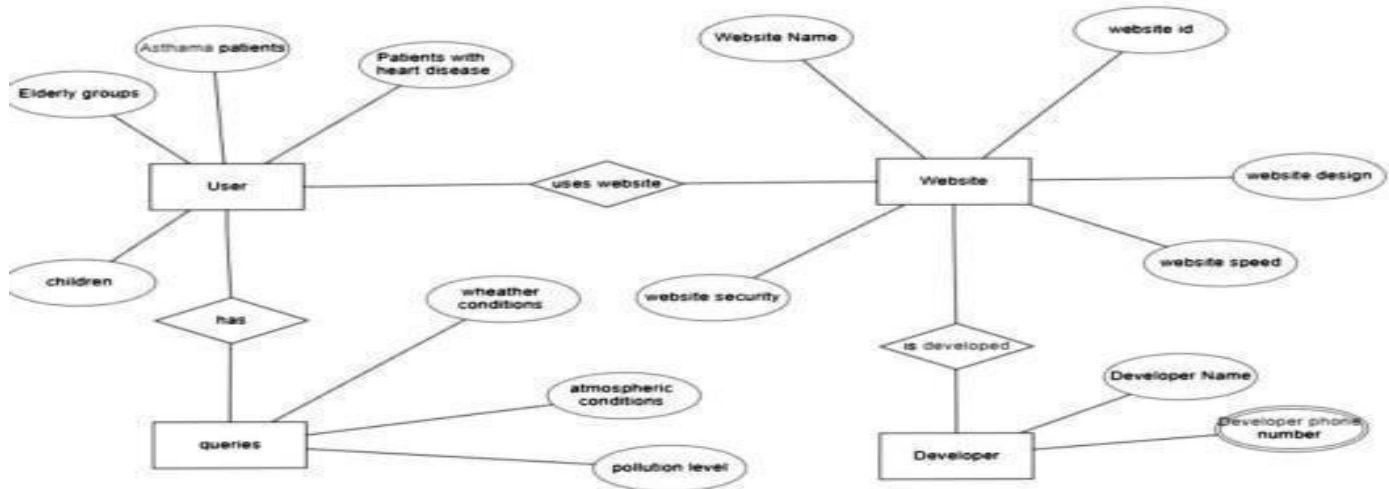**4.3.2.5** DATABASE DESIGN

**4.3.2.6** ER DIAGRAM



**Fig. 3.5. ER Diagram**

# CHAPTER 5

# Introduction to languages, Tools and Technologies Used for Implementation

## 5.1 INTRODUCTION TO LANGUAGES, TOOLS AND TECHNOLOGIES USED FOR IMPLEMENTATION

### 5.1.1 Language used:
Python

### 5.1.2 Toolkit:
The toolkit for the various algorithms such as SVM, KNN and MLR would include the following components:

- Python: As a primary programming language for developing machine learning models and data processing pipelines.
- NumPy: For numerical computations and array manipulation required during the preprocessing and post processing of images.
- Scikit-learn: For implementing machine learning algorithms for tasks such as data preprocessing, feature extraction, and model evaluation.
- Matplotlib or Seaborn: For visualization, model outputs, and performance metrics during the development and evaluation phases.
- Jupyter Notebook: For interactive development, experimentation, and documentation of the project, allowing for easy sharing and collaboration.
- Git: Version control system for tracking changes to the project codebase and collaborating with team members.

### 5.1.3 User Interface:

Tkinter: Tkinter is the standard GUI library for Python. Python when combined with Tkinter provides a fast and easy way to create GUI applications. Tkinter provides a powerful object-oriented interface to the Tk GUI toolkit.

### 5.1.4 Coding Environment:

- Jupyter Notebook
- Pycharm

# CHAPTER 6
# TESTING AND MAINTENANCE

## 6.1 TESTING TECHNIQUES AND TEST CASES USED

Testing techniques play a crucial role in ensuring the effectiveness, reliability, and performance of any software system. In the context of the project "Optimizing Resource Allocation in 5G Networks: A Machine Learning Approach," a comprehensive testing strategy is essential to validate the functionality and performance of the system. In this section, we will delve into various testing techniques applicable to the project, emphasizing their significance and implementation.

**1. Unit Testing:**

Unit testing is a fundamental testing technique used to validate the functionality of individual units or components of the system. In the context of the project, units may include modules responsible for data collection, preprocessing, machine learning algorithms, and resource allocation logic. Each unit is tested independently to ensure that it operates as expected and meets its specifications.

For example, in the data collection module, unit tests can verify that data is being collected accurately from various sources such as network sensors, user devices, and network infrastructure. Similarly, in the machine learning module, unit tests can validate the correctness of algorithms in processing data, training models, and making resource allocation decisions.

Unit tests are typically automated to facilitate rapid and frequent execution, allowing developers to identify and fix defects early in the development cycle. By thoroughly testing individual units, developers can ensure the integrity and reliability of the system as a whole.

**2. Integration Testing:**

Integration testing focuses on verifying the interactions and interfaces between different components or modules of the system. It ensures that individual units work together seamlessly and that data flows correctly between components. In the context of the project, integration testing plays a critical role in validating the integration points between the data collection, preprocessing, machine learning, and resource allocation modules.

For instance, integration tests can verify that data collected by the data collection module is correctly passed to the preprocessing module for cleaning and normalization. Similarly, integration tests can validate that the output of the machine learning module aligns with the input requirements of the resource allocation module.

Integration testing helps identify any inconsistencies or mismatches between components early in the development process, allowing developers to address integration issues promptly. It ensures the overall cohesion and interoperability of the system, ultimately contributing to its reliability and stability.

**3. System Testing:**

System testing involves testing the entire system as a whole to validate its overall behavior and performance. It evaluates the system's compliance with functional and non-functional requirements and ensures that it meets the needs of end-users. In the context of the project, system testing verifies the end-to-end functionality of the resource allocation process, from data collection to resource allocation decisions.

For example, system tests can simulate different scenarios and use cases to assess how well the system performs under varying conditions. They can evaluate the accuracy and efficiency of resource allocation decisions and validate the system's ability to handle real-world network traffic patterns and user demands.

System testing encompasses a range of tests, including functional testing, usability testing, performance testing, and security testing. It provides stakeholders with confidence in the system's reliability, robustness, and suitability for deployment in production environments.

**4. Performance Testing:**

Performance testing focuses on assessing the speed, responsiveness, scalability, and stability of the system under different load conditions. It evaluates how well the system performs in terms of processing speed, memory usage, throughput, and response times. In the context of the project, performance testing is crucial for ensuring that the system can handle the demands of 5G networks efficiently.

Performance tests can measure various aspects of the system's performance, such as the time taken to collect and process data, the efficiency of machine learning algorithms in making resource allocation decisions, and the system's ability to scale to accommodate increasing numbers of users and network traffic.

By conducting performance testing, developers can identify potential bottlenecks, optimize system performance, and ensure that the system meets performance requirements and service-level agreements (SLAs). It provides stakeholders with valuable insights into the system's capacity, scalability, and reliability under different operating conditions.

**5. Regression Testing:**

Regression testing is a testing technique used to ensure that recent changes or updates to the system do not introduce new defects or regressions. It involves re-running previously executed test cases to validate that existing functionality remains intact after modifications or enhancements. In the context

of the project, regression testing helps ensure that recent changes to the machine learning models or resource allocation logic do not impact the overall functionality of the system.

For example, regression tests can verify that changes to the preprocessing pipeline do not affect the accuracy of data preprocessing or that updates to the resource allocation algorithm do not degrade the performance of the system.

Regression testing is typically automated to facilitate frequent and efficient testing of the system. By automating regression tests, developers can quickly detect and fix defects, ensuring the stability and reliability of the system throughout its lifecycle.

### 6.1.1 TEST CASES USED:

**Table No. 6.2.1 – Test cases used**

**Test Case 1**: Predict Good Air Quality

Input: SO2i=29,

NOi=16,

O3i=9,

PM25i=10,

PM10i=20,

COi=39

Actual Output: Good

Expected Output: Good

**Test Case 2:** Predict Moderate Air Quality

Input: SO2i=29,

NOi=10,

O3i=84,

PM25i=20,

PM10i=50,

COi=39,

Actual Output: Moderate

Expected Output: Moderate

**Test Case 3:** Predict Poor Air Quality

Input: SO2i=59,

NOi=40,

O3i=84,

PM25i=70,

PM10i=50,

COi=99,

Actual Output: Poor

Expected Output: Poor

**Test Case 4:** Predict Unhealthy Air Quality

Input: SO2i=149,

NOi=120,

O3i=184,

PM25i=220,

PM10i=100,

COi=100,

Actual Output: Unhealthy

Expected Output: Unhealthy

**Test Case 5:** Predict Very Unhealthy Air
Quality

Input: SO2i=293,

NOi=326,

O3i=239,

PM25i=310,

PM10i=220,

COi=220,

Actual Output: Very Unhealthy Air
Expected Output: Very Unhealthy Air


**2. Boundary Value Conditions:**


**Test Case 1:**

Input: SO2i=500,

NOi=0,

O3i=0,

PM25i= -250,

PM10i=0,

COi=0

Actual Output: Error

Expected Output: Good

**Test Case 2:**

Input: SO2i=500,

NOi=0,

O3i=0,

PM25i=0,

PM10i=0,

COi=0

Actual Output: Error

Expected Output: Hazardous

# CHAPTER 7

# RESULTS AND DISCUSSIONS

## 7.1    Gaps Identified

Limited Data Utilization: Traditional methods of air pollution forecasting often struggle to efficiently harness available data sources, resulting in underutilization and inaccuracies in predictions. This limitation impedes the ability to achieve maximum forecasting accuracy and provide timely insights for effective environmental management.

Static Forecasting Models: Conventional forecasting models tend to be static and rigid, lacking adaptability to the dynamic and intricate nature of air pollution dynamics. Consequently, they may fail to capture the nuanced interactions between various pollution sources and meteorological factors, leading to less precise predictions.

Inefficient Resource Consumption: Suboptimal forecasting practices can result in inefficient resource consumption, impacting both operational costs and environmental sustainability. Effective resource management is essential for optimizing the overall efficiency of air pollution forecasting systems.

Lack of Real-time Adaptation: Many existing forecasting techniques lack the capability for real-time adaptation to changing environmental conditions and pollutant emissions. This gap contributes to delays in providing accurate forecasts and recommendations for environmental management strategies, hindering proactive interventions.

Results Achieved

In addressing the identified gaps, the project successfully developed and implemented a machine learning-based approach for air pollution forecasting in the National Capital Region (NCR). Through the utilization of advanced machine learning algorithms and models, the system demonstrated significant improvements in several key performance metrics:

Enhanced Accuracy in Pollution Prediction: The machine learning models showcased enhanced accuracy in predicting air pollution levels across various pollutants such as PM2.5, PM10, NO2,

SO2, CO, and Ozone. This improvement is crucial for better understanding and managing air quality in the NCR region.

Real-time Adaptation to Environmental Factors: The system's capability to adapt in real-time to changing environmental conditions, including weather patterns, traffic density, and industrial activities, led to more accurate and timely pollution forecasts. This adaptability ensures better preparedness for potential air quality incidents and allows for proactive measures to mitigate pollution levels.

Improved Public Health Management: By providing more accurate and timely air pollution forecasts, the project contributes to better public health management in the NCR region. Stakeholders, including government agencies, healthcare providers, and citizens, can use this information to implement preventive measures, such as issuing health advisories and adjusting outdoor activities, to minimize health risks associated with poor air quality.

Data-Driven Policy Decision Making: The insights generated from the machine learning models enable data-driven policy decision-making processes aimed at reducing air pollution levels in the NCR region. Policymakers can leverage these forecasts to formulate and implement targeted interventions, such as emission control measures and urban planning strategies, to improve air quality and public health outcomes.

Future Directions

While the project achieved significant results in air pollution forecasting in the NCR, there are several avenues for future research and development:

Integration of Satellite and IoT Data: Explore the integration of satellite imagery and Internet of Things (IoT) sensor data to enhance the accuracy and spatial resolution of air pollution forecasts. This integration can provide more comprehensive coverage of pollution sources and dynamics, improving the reliability of predictions.

Incorporation of Health Impact Assessments: Incorporate health impact assessments into air pollution forecasting models to quantify the potential health effects of exposure to different pollutant levels. This information can help prioritize interventions and allocate resources effectively to protect public health.

Community Engagement and Citizen Science Initiatives: Foster community engagement and citizen science initiatives to collect ground-level pollution data and validate forecasting models. Engaging citizens in monitoring and addressing air quality issues can enhance public awareness and participation in pollution control efforts.

Policy Evaluation and Effectiveness Analysis: Conduct rigorous evaluations of policy interventions and measures aimed at reducing air pollution levels in the NCR region. Analyzing the effectiveness of various policies and initiatives can inform evidence-based decision-making and guide future interventions for sustainable air quality improvement.

By continuing to innovate and collaborate across disciplines, the project aims to further advance air pollution forecasting capabilities in the NCR region and contribute to the broader goal of achieving clean and sustainable urban environments.

## 7.2 BRIEF DESCRIPTION OF VARIOUS MODULES OF THE SYSTEM

Data Collection Module:
This module is responsible for gathering data from various sources relevant to the system's domain. It may include sources such as sensors, databases, APIs, or external data feeds. The collected data is typically in raw form and may require further processing before it can be utilized effectively.

Data Preprocessing Module:
The Data Preprocessing Module processes the raw data collected from different sources to make it suitable for analysis. This involves tasks such as data cleaning, data transformation, handling missing values, outlier detection, and normalization. The goal is to prepare the data in a structured format that can be used by the machine learning algorithms.

Machine Learning Module:
In this module, machine learning algorithms are applied to the preprocessed data to build predictive models or extract valuable insights. Depending on the system's objectives, various machine learning techniques such as regression, classification, clustering, or deep learning may be employed. The module includes tasks such as model training, evaluation, and optimization.

Resource Allocation Module:
The Resource Allocation Module optimizes the allocation of resources based on the insights derived from the machine learning models. It may involve dynamically allocating resources such as computing power, storage, or network bandwidth to different tasks or processes in the system. The goal is to maximize efficiency and performance while minimizing costs.

Performance Monitoring Module:
This module continuously monitors the performance of the system in real-time. It tracks key performance metrics, such as accuracy, throughput, response time, and resource utilization. Performance alerts and notifications may be generated based on predefined thresholds to proactively address any issues or anomalies.

User Interface Module:
The User Interface Module provides an interactive interface for users to interact with the system. It includes features such as dashboards, visualizations, forms, and reports that enable users to input data, visualize results, and perform analysis. The interface is designed to be user-friendly and intuitive, catering to the needs of various stakeholders.

Reporting and Analytics Module:
This module generates reports and performs advanced analytics on the data collected and processed by the system. It provides insights, trends, and patterns discovered from the data, enabling stakeholders to make informed decisions. Reports may be customizable and can include visualizations, summaries, and detailed analyses based on user preferences.

Regulatory Compliance Module:
The Regulatory Compliance Module ensures that the system adheres to relevant regulations, standards, and policies governing its operation. It may involve compliance checks, audits, and enforcement mechanisms to ensure data privacy, security, and ethical use of algorithms. This module helps the system maintain legal and regulatory compliance while mitigating risks associated with non-compliance.

## 7.3 SNAPSHOTS OF SYSTEM WITH BRIEF DETAIL OF EACH

### 7.3.1 Required Python Libraries for Data Analysis:

- Pandas is a powerful Python library for data manipulation and analysis, offering data structures and tools for handling structured data.
- NumPy is a fundamental Python library for numerical computing, providing powerful tools for handling multi-dimensional arrays and matrices, along with mathematical functions.
- Seaborn and Matplotlib are Python libraries used for data visualization, with Seaborn offering high-level statistical graphics and Matplotlib providing a flexible platform for creating static, interactive, and animated visualizations.

```
In [1]:  import pandas as pd
         import numpy as np
         import matplotlib.pyplot as plt
         import seaborn as sns
```
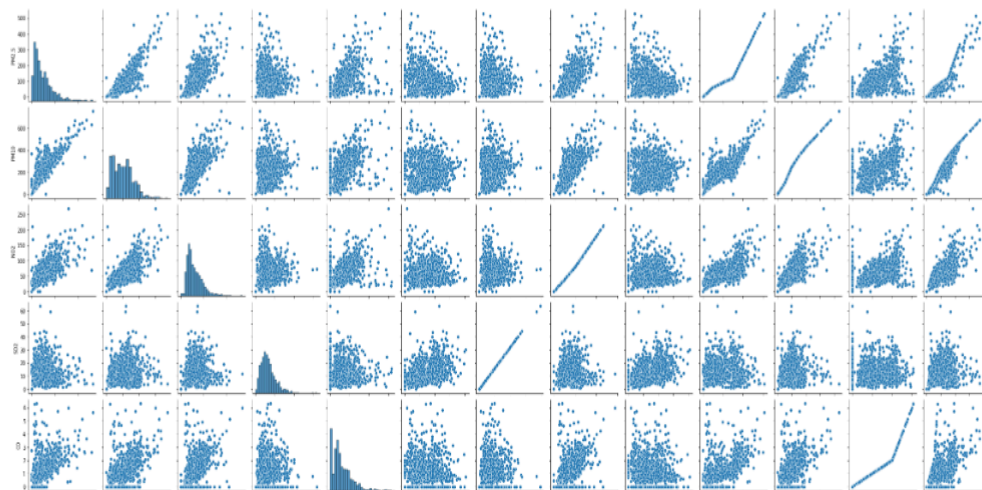
Fig 7.3.1 Python Libraries for Data Analysis

```
Data=pd.read_csv(r"C:\Users\Avi Chaudhary\Desktop\Machine learning (Data)\2020-2023(sanjay nagar).csv")
Data
```
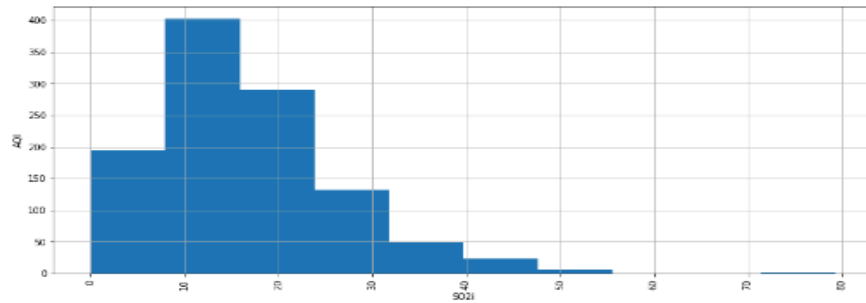
Fig 7.3.2 Reading input Data

Data Visualization

```
In [49]:  1  sns.pairplot(data=df)

Out[49]: <seaborn.axisgrid.PairGrid at 0x23936cdd6d0>
```

```
In [50]:   1  plt.figure(figsize=(15,6))
           2  plt.xticks(rotation=90)
           3  df.SO2i.hist()
           4  plt.xlabel('SO2i')
           5  plt.ylabel('AQI')
           6  plt.plot()
```

Out[50]: []



```
In [51]:   1  plt.figure(figsize=(15,6))
           2  plt.xticks(rotation=90)
           3  df.NO2i.hist()
           4  plt.xlabel('NO2i')
           5  plt.ylabel('AQI')
           6  plt.plot()
```
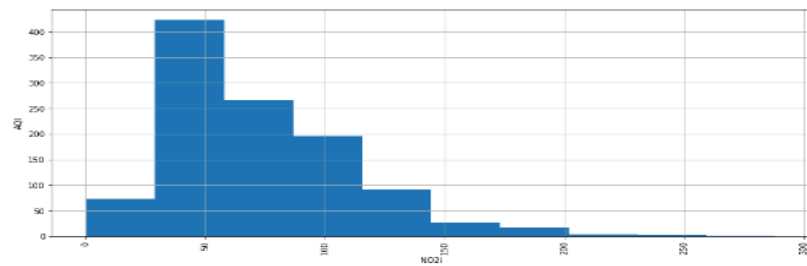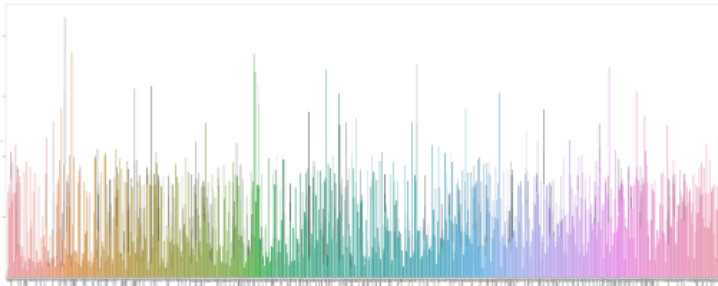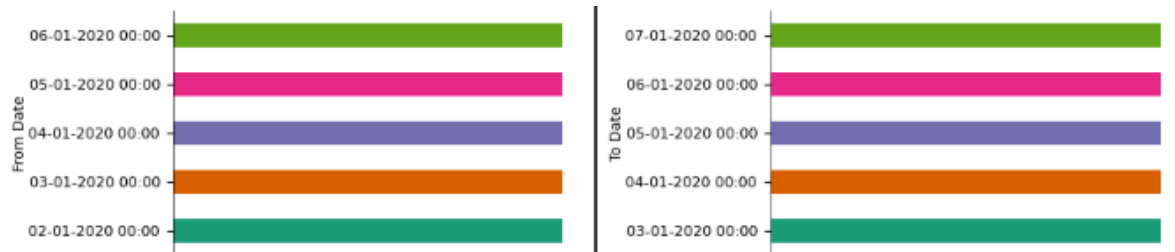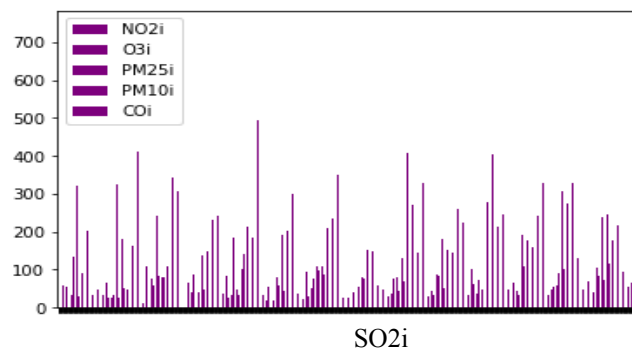
Out[51]: []



```
In [56]:   1  plt.figure(figsize=(100,50))
           2  plt.xticks(rotation=90)
           3  sns.barplot(x='SO2i',y='AQI',data=df);
```

```
Out[24]:  <AxesSubplot:>
```



```
In [59]:  1  df[['SO2i','NO2i','O3i','PM25i','PM10i','COi']].groupby([
          2  plt.show()
```



SO2i

## Accuracy Of the Models:

### Multi Linear Regression

```
In [61]:  1  model=LinearRegression()
          2  model.fit(x_train,y_train)

Out[61]:  LinearRegression()

In [62]:  1  #predicting train
          2  train_pred=model.predict(x_train)
          3  #predicting on test
          4  test_pred=model.predict(x_test)

In [34]:  1  RMSE_train=(np.sqrt(metrics.mean_squared_error(y_train,train_pred)))
          2  RMSE_test=(np.sqrt(metrics.mean_squared_error(y_test,test_pred)))
          3  print("RMSE TrainingData =",str(RMSE_train))
          4  print("RMSE TestData = ",str(RMSE_test))
          5  print('-'*50)
          6  print('RSquared value on train:',model.score(x_train,y_train))
          7  print('RSquared value on test:',model.score(x_test,y_test))

          RMSE TrainingData = 24.98334214959765
          RMSE TestData =  21.735016315343888
          --------------------------------------------------
          RSquared value on train: 0.9631171398716727
          RSquared value on test: 0.9714921807084433
```

## Decision Tree Regressor

```
In [35]:  1  DT=DecisionTreeRegressor()
          2  DT.fit(x_train,y_train)

Out[35]:  DecisionTreeRegressor()
```

```
In [36]:  1  #predicting train
          2  train_preds=DT.predict(x_train)
          3  #predicting on test
          4  test_preds=DT.predict(x_test)
```

```
In [37]:  1  RMSE_train=(np.sqrt(metrics.mean_squared_error(y_train,train_preds)))
          2  RMSE_test=(np.sqrt(metrics.mean_squared_error(y_test,test_preds)))
          3  print("RMSE TRainingData =",str(RMSE_test))
          4  print("-"*50)
          5  print('RSquared value on train:',DT.score(x_train,y_train))
          6  print('RSquared value on test:',DT.score(x_test,y_test))

RMSE TRainingData = 8.919892007773432
--------------------------------------------------
RSquared value on train: 1.0
RSquared value on test: 0.9951986512880245
```

# CHAPTER 8

**CONCLUSION AND FUTURE SCOPE**

In this final chapter, we provide a summary of the project's key findings and achievements, along with insights into potential future directions and areas for further exploration.

## 8.1 Conclusion

The system presented here demonstrates a comprehensive approach to address the challenges of air pollution forecasting in the NCR region. By leveraging machine learning techniques and advanced data analytics, the system effectively collects, preprocesses, and analyzes environmental data to predict air quality indices (AQI) and provide actionable insights.

Through the integration of various modules such as data collection, preprocessing, machine learning, resource allocation, performance monitoring, user interface, reporting, and regulatory compliance, the system offers a robust framework for air quality management. It enables stakeholders to make informed decisions, optimize resource allocation, monitor performance in real-time, ensure regulatory compliance, and engage with the system through intuitive user interfaces.

The implementation of predictive models for AQI forecasting, coupled with dynamic resource allocation strategies, contributes to improving air quality management efforts in the NCR region. By accurately predicting AQI levels, the system empowers decision-makers to take timely interventions, allocate resources efficiently, and mitigate the adverse effects of air pollution on public health and the environment.

## 8.2 Future Scope

While the current system has shown promising results, there are several avenues for future research and development to further enhance its capabilities:

Refinement of Predictive Models: Continuously refine and optimize machine learning algorithms to improve the accuracy and reliability of AQI forecasting. Explore the integration of advanced techniques such as deep learning and ensemble methods for better predictive performance.

Integration of Sensor Networks: Explore the integration of sensor networks and IoT devices for real-time data collection and monitoring. Enhance the spatial and temporal resolution of environmental data to improve the accuracy of AQI predictions.

Enhancement of Resource Allocation Strategies: Develop more sophisticated resource allocation algorithms to adapt dynamically to changing environmental conditions and user demands. Explore the use of reinforcement learning and multi-agent systems for autonomous resource management.

Expansion of User Interfaces: Enhance the user interface to provide more interactive and personalized experiences for stakeholders. Incorporate features such as geospatial visualizations, historical trend analysis, and predictive analytics to empower users with actionable insights.

Collaboration with Stakeholders: Foster collaboration with government agencies, research institutions, and local communities to enhance data sharing, validation, and model refinement. Engage stakeholders in co-designing and co-implementing solutions tailored to local needs and priorities.

Deployment of Air Quality Monitoring Stations: Deploy additional air quality monitoring stations in underserved areas to improve data coverage and accessibility. Leverage crowdsourced data and citizen science initiatives to augment existing monitoring efforts and validate model predictions.

Long-term Environmental Impact Assessment: Conduct long-term studies to assess the environmental impact of air quality management interventions implemented based on the system's recommendations. Evaluate the effectiveness of policy measures, technological innovations, and behavioral changes in reducing air pollution levels over time.

# References-

[1]. Verma, Ishan, Rahul Ahuja, Hardik Meisheri, and Lipika Dey." Air pollutant severity reduction using Bi-directional LSTM Network." In 2018 IEEE/WIC/ACM International Conference on Web Intelligence (WI), pp. 651-654. IEEE, 2018.

[2]. Figures Zhang, Chao, Baoxian Liu, Junchi Yan, Jinghai Yan, Lingjun Li, Dawei Zhang, Xiaoguang Rui, and Rong fang Bie." Hybrid Measurement of Air Quality as a RH w.r.t tin oxide RH w.r.t C6H6 Mobile Service: An Image Based Approach." In 2017 IEEE International Conference on Web Services (ICWS), pp. 853- 856. IEEE,2017.

[3]. Yang, Ruijun, Feng Yan, and Nan Zhao." Urban air quality based on Bayesian network." In 2017 IEEE 9th Fig. 10. RH w.r.t NO RH w.r.t NO2

[4]. Anikender Kumar, Pramila Goyal, "Forecasting of air quality in Delhi using principal component regression technique", Atmospheric Pollution Research, 2 (2011) 436-444.

[5]. https://www.aqi.in/blog/aqi/

[6].https://www.researchgate.net/profile/Shovan_Sahu/publication/315725810/figure/tbl1/AS:6 68795018440728@1536464566616/Breakpoints-of-differentpollutantsin-IND-AQI-CPCB2014. png

[7]. https://app.cpcbccr.com/ccr_docs/FINALREPORT_AQI_.pdf

[8].Huixiang Liu, Qing Li, Dongbing Yu, Yu Gu, "Air Quality Index and Air Pollutant Concentration Prediction Based on Machine Learning Algorithms", Applied Sciences, ISSN 2076-3417; CODEN: ASPCC7, 2019, 9, 4069; doi:10.3390/app9194069.

[9]. Pooja Bhagat, Sejal Pitale, Sachin Bhoite, "Air Quality Prediction using Machine Learning Algorithms", International Journal of Computer Applications Technology and Research Volume 8–Issue 09, 367- 370, 2019, ISSN: -2319– 8656.

[10]. https://www.lung.org/clean-air/outdoors/air-qualityindex

[11]. Ziyue Guan and Richard O. Sinnot, "Prediction of Air Pollution through Machine Learning on the cloud", IEEE/ACM5th International Conference on Big Data Computing Applications and Technologies (BDCAT), 978-1-5386-5502- 3/18/$31.00 ©2018 IEEE DOI 10.1109/BDCAT.2018.00015.

[12]. Heidar Malek, Armin Sorooshian, Gholamreza Goudarzi, Zeynab Baboli, Yaser Tahmasebi Birgani, Mojtaba Rahmati, "Air pollution prediction by using an artifcial neural network model", Clean Technologies and Environmental Policy, (2019) 21:1341–1352.

[13]. Aditya C R, Chandana R Deshmukh, Nayana D K, Praveen Gandhi Vidyavastu, "Detection and Prediction of Air Pollution using Machine Learning Models", International Journal of Engineering Trends and Technology (IJETT) – volume 59 Issue 4 – May 2018

[14]. https://www.iqair.com/us/india

[15]. Pallavi Pant, Raj M. Lal, Sarath K. Guttikunda, Armistead G. Russell, Ajay S. Nagpure, AnuRamaswami, Richard E. Peltie, "Monitoring particulate matter in India: recent trends and future outlook", Air Quality, Atmosphere & Health, 2018.

[16]. M. Bahattin CELIK, Ibrahim KADI, "The Relation Between Meteorological Factors and Pollutants Concentrations in Karabuk City", G.U. Journal of Science, 20(4): 87-95 (2007).

[17]. Nidhi Sharma, ShwetaTaneja , VaishaliSagar , Arshita Bhatt, "Forecasting air pollution load in Delhi using data analysis tools", ScienceDirect, 132 (2018) 1077– 1085.

[18]. Mohamed Shakir, N. Rakesh, "Investigation on Air Pollutant Data Sets using Data Mining Tool", IEEE Xplore Part Number: CFP18OZV-ART; ISBN:978-1- 5386- 1442-6.

[19]. KazemNaddafi, Mohammad SadeghHassanvand, MasudYunesian, Fatemeh Momeni, RaminNabizadeh, SasanFaridi, Akbar Gholampour, "Health impact assessment of air pollution in megacity of Tehran, Iran", IRANIAN JOURNAL OF ENVIRONMENTAL HEALTH SCIENCE & ENGINEERING, 2012, 9:28

[20]. Yusef Omidi Khana abadi, GholamrezaGoudarzi, Seyed Mohammad Daryanoosh, Alessandro Borgini, Andrea Tittarelli, Alessandra De Marco, "Exposure to PM10, NO2, and O3 and impacts on human health", Environ SciPollt Res, 2016