

Machine Learning Approaches to Assess the Health Impact of Air Pollution

BY

Md. Rahat Mahabub Anik

ID: 191-25-730

This Report Presented in Partial Fulfillment of the Requirements for
The Degree of Masters of Science in Computer Science and Engineering

Supervised By

Dr. Sheak Rashed Haider Noori

Professor & Head.

Department of CSE

Daffodil International University

Co Supervised By

Abdus Sattar

Assistant Professor & Coordinator MSc.

Department of CSE

Daffodil International University



DAFFODIL INTERNATIONAL UNIVERSITY

DHAKA, BANGLADESH

JANUARY 2025

APPROVAL

This Project/Thesis titled “Machine Learning Approaches to Assess the Health Impact of Air Pollution”, submitted by Md. Rahat Mahabub Anik, ID No: 191-25-730 to the Department of Computer Science and Engineering, Daffodil International University has been accepted as satisfactory for the partial fulfillment of the requirements for the degree of M.Sc. in Computer Science and Engineering and approved as to its style and contents. The presentation has been held on 11-01-2025.

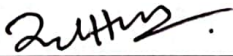
BOARD OF EXAMINERS



Chairman

Dr. Sheak Rashed Haider Noori, PhD
Professor and Head
Department of Computer Science and Engineering
Faculty of Science & Information Technology
Daffodil International University

Internal Examiner



Dr. Md. Zahid Hasan, PhD
Associate Professor
Department of Computer Science and Engineering
Faculty of Science & Information Technology
Daffodil International University

Internal Examiner



Dr. Arif Mahmud, PhD
Associate Professor & Director MIS
Department of Computer Science and Engineering
Faculty of Science & Information Technology
Daffodil International University

External Examiner



Dr. Mohammed Nasir Uddin, PhD
Professor
Department of Computer Science and Engineering
Jagannath University

DECLARATION

We hereby declare that, this thesis has been done by us under the supervision of **Dr. Sheak Rashed Haider Noori, Professor and Head, Department of CSE Daffodil International University**. We also declare that neither this thesis nor any part of this thesis has been submitted elsewhere for award of any degree or diploma.

Supervised by:



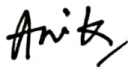
Dr. Sheak Rashed Haider Noori
Professor & Head
Department of CSE
Daffodil International University

Co Supervised by:



Abdus Sattar
Assistant Professor & Coordinator MSc.
Department of CSE
Daffodil International University

Submitted by:



Md. Rahat Mahabub Anik
ID: 191-25-730
Department of CSE
Daffodil International University

ACKNOWLEDGEMENT

First, we express our heartiest thanks and gratefulness to almighty God for His divine blessing makes us possible to complete the final thesis successfully.

We really grateful and wish our profound our indebtedness to **Dr. Sheak Rashed Haider Noori, Professor & Head**, Department of CSE Daffodil International University, Dhaka. Deep Knowledge & keen interest of our supervisor in the field of “Data Mining” to carry out this thesis. His endless patience, scholarly guidance, continual encouragement, constant and energetic supervision, constructive criticism, valuable advice, reading many inferior draft and correcting them at all stage have made it possible to complete this thesis.

We would like to express our heartiest gratitude to **Dr. Sheak Rashed Haider Noori, Professor and Head, Department of CSE**, for his kind help to finish our thesis and also to other faculty members and the staff of CSE department of Daffodil International University.

We would like to thank our entire course mate in Daffodil International University, who took part in this discuss while completing the course work.

Finally, we must acknowledge with due respect the constant support and passion of our parents.

ABSTRACT

Air pollution is a critical global health concern, linked to a wide range of respiratory, cardiovascular, and chronic diseases. This study investigates the use of machine learning to assess the health impacts of air pollution by leveraging diverse datasets that include air quality indices, meteorological factors, and health outcomes. The research evaluates several algorithms, including Decision Trees, Support Vector Machines, Random Forest, XGBoost, and Tree Ensemble, using performance metrics such as precision, recall, F1 score, and accuracy. Among these, ensemble methods, particularly Tree Ensemble, demonstrated superior generalization capabilities, achieving an accuracy of 97%. The study highlights the ability of machine learning models to capture complex, non-linear relationships between environmental and health variables, offering significant improvements over traditional statistical methods. However, challenges such as data heterogeneity, ethical concerns, and model interpretability remain critical barriers. The research emphasizes the need for incorporating real-time data streams, explainable AI techniques, and fairness mechanisms to enhance model transparency and usability for decision-making. This work underscores the transformative role of machine learning in mitigating the adverse health effects of air pollution. By integrating predictive models with policy frameworks and fostering global collaboration, future studies can drive effective public health interventions and promote sustainable environmental practices.

LIST OF FIGURES

FIGURES	PAGE NO
Figure 1: Sample Dataset	17
Figure 2: Distribution of HealthImpactClass	18
Figure 3: Heatmap of the Dataset	18
Figure 4: Distribution of AQI	19
Figure 5: Distribution of PM10	19
Figure 6: Distribution of PM2_5	20
Figure 7: Distribution of NO2	20
Figure 8: Distribution of SO2	21
Figure 9: Distribution of O3	21
Figure 10: Distribution of Humidity	22
Figure 11: Distribution of RespiratoryCases	22
Figure 12: Distribution of HospitalAdmissions	23
Figure 13: Proposed Methodology	23
Figure 14: Score of Feature Selection	24
Figure 15: Comparison of Training and Testing Accuracy of each Algorithm	30
Figure 16: Performance Comparison of each Algorithm	31
Figure 17: Confusion Matrix for Decision Tree	32
Figure 18: Confusion Matrix for Random Forest	33
Figure 19: Confusion Matrix for KNN	33
Figure 20: Confusion Matrix for Naïve Bayes	34
Figure 21: Confusion Matrix for SVM	34
Figure 22: Confusion Matrix for XGBoost	35
Figure 23: Confusion Matrix for Gradient Boosting	35
Figure 24: Confusion Matrix for Cat Boost	36
Figure 25: Confusion Matrix for Ada Boost	36
Figure 26: Confusion Matrix for MLP Classifier	37
Figure 27: Confusion Matrix for Tree Ensemble	37
Figure 28: ROC Curve for Tree Ensemble	38

LIST OF TABLES

TABLE	PAGE NO
Table 1: Related Work	10
Table 2: Training and Testing Accuracy of Used Algorithms	30
Table 3: Training and Testing Accuracy of Used Algorithms	31

TABLE OF CONTENTS

CONTENTS	PAGE
Approval Page	i
Declaration	ii
Acknowledgements	iii
Abstract	iv
List of Figures	vii
List of Tables	vii
CHAPTER	
CHAPTER 1: INTRODUCTION	1-4
1.1 Introduction	1
1.2 Motivation	1-2
1.3 Research Objective	2
1.4 Research Question	3
1.5 Report Layout	3-4
CHAPTER 2: BACKGROUND	5-14
2.1 Introduction	5
2.2 Related Work	5-10
2.3 Scope of the problem	11-12
2.4 Challenges	12-14
CHAPTER 3: RESEARCH METHODOLOGY	15-28
3.1 Introduction	15
3.2 Research Subject and Instrumentation	16-17
3.3 Data Collection	17
3.4 Statistical Analysis	17-23
3.5 Proposed Methodology	23-28

3.6 Implementation Requirements	28
CHAPTER 4: EXPERIMENTAL RESULTS AND DISCUSSION	29-39
4.1 Introduction	29
4.2 Experimental Results	29-38
4.3 Discussion	38-39
CHAPTER 5: IMPACT ON SOCIETY, ENVIRONMENT, AND SUSTAINABILITY	40-45
5.1 Impact on Society	40-41
5.2 Impact on Environment	41-42
5.3 Ethical Aspects	42-43
5.4 Sustainability Plan	43-45
CHAPTER 6: CONCLUSION AND FUTURE WORK	46-49
6.1 Summary of the study	46
6.2 Conclusions	46-47
6.3 Implication for further study	47-48
REFERENCES	49-50

CHAPTER 1

INTRODUCTION

1.1 Introduction

The introduction begins by emphasizing the growing global concern about air pollution and its adverse health effects, ranging from respiratory diseases to cardiovascular complications. It highlights the importance of accurate and timely assessment of air pollution's impact on public health to inform policies and interventions. Traditional epidemiological approaches, while valuable, often face challenges such as limited data availability, high costs, and the inability to capture real-time or granular variations in air quality and health outcomes.

Machine learning (ML) emerges as a transformative tool in this context, offering the capability to process vast datasets, identify complex patterns, and predict health outcomes with high precision. The section underscores how ML techniques can integrate diverse data sources, such as satellite imagery, sensor networks, and healthcare records, to provide a more holistic understanding of the air pollution-health relationship [1].

The introduction also outlines the scope of the paper, which focuses on exploring various ML approaches—such as supervised learning, unsupervised learning, and deep learning—applied to the assessment of air pollution's health impact. It concludes by discussing the potential of ML-driven insights to guide effective public health strategies and address gaps in traditional assessment methods, setting the stage for the detailed exploration of methodologies, datasets, and results in subsequent sections [3].

1.2 Motivation

The escalating levels of air pollution have become a significant global health concern, contributing to millions of premature deaths and a growing burden of respiratory, cardiovascular, and other chronic diseases [2]. Traditional methods for assessing the health impact of air pollution, while insightful, are often constrained by limitations in data availability, scalability, and the ability to capture complex interactions between environmental and health variables. This underscores the need for innovative approaches to improve the precision and efficiency of impact assessments.

Machine learning (ML) presents a compelling solution to these challenges. With its ability to analyze large, diverse datasets, ML can uncover patterns and correlations that are not easily identifiable using conventional techniques. The integration of data from various sources—such

as air quality monitoring systems, satellite imagery, and healthcare records—provides an opportunity to develop predictive models capable of real-time and localized assessments.

Furthermore, the growing availability of open-source tools, computational resources, and public datasets has lowered the barriers to implementing ML in this domain [4]. This study is motivated by the potential of ML to enhance our understanding of the health impacts of air pollution, thereby supporting the development of targeted public health policies and interventions. By exploring and evaluating ML approaches, this work aims to bridge the gap between traditional epidemiological methods and advanced computational techniques, ultimately contributing to more effective strategies for mitigating air pollution's health consequences.

1.3 Research Objective

The main objective of our project to develop a reliable machine learning based model to assessing the health impact of air pollution.

The objectives of this research are:

- To aggregate and preprocess diverse datasets, including air quality monitoring data, satellite imagery, meteorological records, and healthcare datasets, for comprehensive analysis.
- To analyze the complex, nonlinear relationships between pollutant levels and health impacts, identifying key pollutants and risk factors influencing public health.
- To develop and evaluate machine learning models (e.g., supervised learning, unsupervised learning, deep learning) for accurately predicting health outcomes associated with air pollution exposure.
- To provide actionable insights and recommendations to support evidence-based public health policies and interventions aimed at mitigating the health impacts of air pollution.

By achieving these objectives, the study aims to enhance the understanding of air pollution's health impacts and support data-driven solutions for global and regional public health challenges.

1.4 Research Questions

Based on the objectives you've provided for "Machine Learning Approaches to Assessing the Health Impact of Air Pollution.", here are some potential research questions that could guide your investigation:

RQ1. How can machine learning techniques effectively integrate diverse datasets (e.g., air quality data, satellite imagery, healthcare records) to assess the health impact of air pollution?

RQ2. Can machine learning approaches capture complex, nonlinear relationships between pollutant concentrations and various health outcomes?

RQ3. To what extent can machine learning models provide real-time predictions of health risks based on current air pollution levels?

RQ4. Which machine learning models (e.g., supervised, unsupervised, deep learning) are most effective in predicting health outcomes associated with air pollution exposure?

1.5 Report Layout

The report has total 6 Chapters which will be followed given by instructions:

In Chapter 1, we have discussed about the introduction of this research in this part of this research, we discussed in details about the importance of sales forecasting. Research question, motivation of this research, rationale and expected output has also discussed in this part.

In Chapter 2, we reviewed existing work on leaf disease prediction in this chapter. We discussed about other authors approaches, limitations, results, methods in this part. Research scope and challenges also have discussed here.

In Chapter 3, research methodology has mainly discussed here. This chapter shows the data collection procedure, statistics of data, classifiers, figures of stats. Implementation requirements also have discussed in this chapter.

In Chapter 4, results of this research have mainly showed. In this chapter we have showed the result of all the classifiers we have used in this research.

In Chapter 5, shows how this research can impact in our society. Why this work is important and how it will sustain in this arena, this chapter shows mainly.

In Chapter 6, discussion and conclusion of this research have shown. In this part we have also discussed about the future work, and the limitations of this work.

CHAPTER 2

BACKGROUND

2.1 Introduction

The literature review provides a comprehensive examination of existing studies that employ machine learning techniques to evaluate the health impacts of air pollution. It begins by summarizing traditional approaches to air pollution impact assessment, highlighting their limitations in handling large-scale, dynamic datasets and capturing complex relationships between variables. This establishes the need for machine learning as a complementary or alternative methodology.

The review then explores the application of supervised learning models, such as regression and classification algorithms, which have been widely used to predict specific health outcomes based on pollutant levels. It discusses advancements in unsupervised learning and clustering techniques for uncovering hidden patterns in air quality and health data. Deep learning models, particularly convolutional neural networks (CNNs) and recurrent neural networks (RNNs), are also analyzed for their capability to process spatial, temporal, and high-dimensional datasets, such as satellite images and time-series pollution data.

Key themes in the literature include the integration of diverse datasets (e.g., satellite imagery, air quality sensors, healthcare records), the development of real-time prediction systems, and efforts to improve the interpretability and scalability of machine learning models. The review identifies gaps in the existing research, such as limited focus on regional and demographic variations, challenges in model generalization, and the need for interdisciplinary collaboration.

By synthesizing the findings from prior studies, the literature review sets the stage for this research to address identified gaps and push the boundaries of how machine learning can enhance our understanding of the health impacts of air pollution.

2.2 Related Works

The increasing prevalence of air pollution and its adverse effects on public health have motivated significant research into understanding the complex relationships between environmental exposure and health outcomes. Traditional epidemiological methods, while foundational, often face limitations in processing large, diverse datasets and capturing nonlinear relationships. Machine learning (ML) has emerged as a transformative tool in this

domain, enabling researchers to analyze complex datasets, identify patterns, and make accurate predictions.

With the rapid pace of urbanization, many developing countries are facing significant challenges due to severe air pollution. Accurately predicting future air quality has become increasingly crucial for both governmental policy-making and individual decision-making. This paper [1] focuses on forecasting air quality for the next 48 hours at each monitoring station, utilizing air quality, meteorological, and weather forecast data. Leveraging domain knowledge in air pollution, we introduce DeepAir, a deep neural network (DNN)-based model. DeepAir integrates a spatial transformation module and a deep distributed fusion network. The spatial transformation module accounts for the spatial correlations of air pollutants, converting sparse air quality data into a uniform input to simulate pollution sources. The deep distributed fusion network, on the other hand, employs a neural distributed architecture to combine diverse urban data, such as meteorological conditions, capturing the key factors influencing air quality. DeepAir has been implemented in the Air Pollution Prediction system, delivering precise air quality forecasts for over 300 cities in China on an hourly basis. Experimental evaluations using three years of data from nine Chinese cities highlight DeepAir's superiority over 10 baseline methods. Compared to the prior online approach in the Air Pollution Prediction system, DeepAir achieves relative accuracy improvements of 2.4% for short-term, 12.2% for long-term, and 63.2% for sudden change predictions.

In the study [5] the authors evaluated the effectiveness of Beijing's 2013–2017 Clean Air Action Plan, which aimed to reduce air pollutant emissions and improve ambient air quality. Traditional assessment methods, such as statistical and chemical transport modeling, often face uncertainties. To address this, the study employed a machine-learning-based random forest technique to separate the effects of meteorological conditions from actual emission reductions. The findings revealed that, after accounting for meteorological influences, there were significant reductions in pollutants: PM_{2.5} decreased by approximately 34%, PM₁₀ by 24%, NO₂ by 17%, SO₂ by 68%, and CO by 33% between 2013 and 2017. The substantial decline in PM_{2.5} and SO₂ was largely attributed to reduced coal combustion. The study concluded that the action plan was highly effective in lowering primary pollution emissions and improving air quality in Beijing, offering a successful model for developing air quality policies in other regions.

The study [6] explores the application of machine learning techniques to forecast air pollutant levels and the Air Quality Index (AQI) in California. Given the dynamic and volatile nature of air pollutants, accurately predicting air quality is complex yet crucial for public health

and environmental protection. The researchers employed Support Vector Regression (SVR) with a Radial Basis Function (RBF) kernel to model and predict concentrations of pollutants such as carbon monoxide, sulfur dioxide, nitrogen dioxide, ground-level ozone, and particulate matter (PM_{2.5}). The SVR model demonstrated high accuracy in forecasting hourly pollutant concentrations and AQI values. Notably, the model achieved a 94.1% accuracy rate in classifying AQI into six categories defined by the U.S. Environmental Protection Agency on unseen validation data. The study also found that utilizing the complete set of available variables yielded better predictive performance than feature selection methods like Principal Component Analysis. These findings underscore the potential of machine learning models in effectively predicting air quality, thereby aiding in timely public health interventions and policy-making.

The study [7] investigates the effectiveness of various machine learning algorithms in forecasting air quality within urban environments. Recognizing the critical role of accurate air quality predictions for public health and environmental management, the researchers evaluated multiple models, including Random Forest, Decision Tree, Linear Regression, and XGBoost. Their analysis revealed that ensemble methods, particularly Random Forest and XGBoost, outperformed other techniques in terms of prediction accuracy and robustness. The study underscores the importance of selecting appropriate machine learning models to enhance air quality monitoring systems in smart cities, thereby facilitating timely interventions and informed policy decisions.

The study [8] introduces a cost-effective, field-portable platform named c-Air, designed for high-throughput quantification of particulate matter (PM) in the air. This innovative device integrates computational lens-free microscopy with machine learning algorithms to rapidly screen air samples and generate microscopic images of aerosols. Capable of analyzing 6.5 liters of air in just 30 seconds, c-Air provides detailed statistics on particle size and density distributions with an accuracy of approximately 93%. The system is also equipped with a smartphone application for device control and real-time display of results. Field tests conducted in various indoor and outdoor environments demonstrated strong correlation between c-Air measurements and those obtained from Environmental Protection Agency-approved devices. Additionally, c-Air was utilized to map air quality around Los Angeles International Airport, revealing significant PM concentration increases even at distances greater than 7 kilometers from the airport, particularly along flight paths. The study highlights c-Air's adaptability in detecting specific airborne particles, such as various types of pollen and mold, offering a promising solution for distributed and accurate air quality monitoring.

In the paper [10] address the challenge of forecasting hourly concentrations of air pollutants such as ozone, PM_{2.5}, and sulfur dioxide. The authors propose a machine learning framework that formulates the prediction of 24-hour pollutant concentrations as a multi-task learning (MTL) problem, allowing for the simultaneous modeling of multiple related tasks. To enhance model performance and prevent overfitting, they introduce a novel regularization technique that enforces similarity between prediction models of consecutive hours. This approach is compared with standard regularization methods, including Frobenius norm, nuclear norm, and $\ell_{2,1}$ -norm regularizations. Experimental results demonstrate that the proposed regularization method, combined with parameter-reducing formulations, outperforms existing regression models and traditional regularization techniques in predicting hourly air pollution concentrations.

This study [11] developed predictive models for air quality using an extensive 11-year dataset from Taiwan's Environmental Protection Administration. The research aimed to address the limitations of previous studies that often relied on shorter datasets, which inadequately captured seasonal and other temporal variations affecting air quality. By employing machine learning techniques, the authors created models capable of forecasting instances of poor air quality with improved accuracy. The study's findings underscore the importance of long-term data in enhancing the reliability of air quality predictions, thereby contributing valuable insights for environmental monitoring and public health planning.

Developing models to evaluate air pollution exposure within urban areas has been recognized as a critical focus for future research, particularly for assigning exposure levels in health studies. This paper [13] reviews six categories of models used for assessing intra-urban air pollution exposure: (i) proximity-based methods, (ii) statistical interpolation techniques, (iii) land-use regression models, (iv) line dispersion models, (v) integrated emission-meteorological models, and (vi) hybrid models that combine personal or household exposure monitoring with one of the aforementioned approaches. The review is enriched with practical examples from Hamilton, Canada, and includes a qualitative evaluation of the models based on key criteria relevant to health effects research. Hybrid models show significant potential for addressing the challenge of achieving population-representative samples while capturing individual-level exposure variations. Furthermore, emerging tools like remote sensing and activity-space analysis are likely to enhance the precision of traditional methods. As advancements in these approaches continue, the field of exposure assessment may play a vital role in reducing the scientific uncertainties currently limiting effective public health policy interventions.

This study [14] presents a causal inference framework for estimating the number of adverse health events prevented by large-scale air quality regulations through reductions in exposure to multiple pollutants. The approach is motivated by regulations that affect pollution levels across all areas under their jurisdiction. We define a new causal estimand, the Total Events Avoided (TEA), as the difference between the expected number of health events under hypothetical no-regulation pollution exposures and the observed number of health events under with-regulation pollution exposures. To estimate TEA, we propose two methods: a matching method and a machine learning method, both utilizing high-resolution, population-level data on pollution and health outcomes. These methods improve on traditional regulatory health impact analyses by explicitly addressing causal assumptions, leveraging population-level data, reducing reliance on parametric assumptions, and accounting for the combined effects of multiple pollutants. To ensure robust and conservative estimates, the TEA calculation focuses on health impacts for units whose anticipated no-regulation features fall within the range of observed with-regulation data. This minimizes model dependence and complements traditional parametric methods with a data-driven perspective. We demonstrate the application of these methods by analyzing the health impacts of the 1990 Clean Air Act Amendments on the US Medicare population, offering a refined approach to evaluating the benefits of air quality regulations.

Air pollution poses a significant global health risk, contributing to respiratory diseases. Advances in air quality mapping, supported by smart city initiatives and the proliferation of IoT sensor devices, have increased data availability and driven progress in air pollution forecasting. This study [15] proposes an integrated approach to predicting air quality using image data and assessing lung disease severity based on the Air Quality Index (AQI). The objective is to refine existing techniques for improved accuracy in forecasting AQI and related pollutants, including PM_{2.5}, PM₁₀, O₃, CO, SO₂, and NO₂, and to evaluate lung disease severity. The approach combines a VGG16 model for feature extraction from images with a neural network for AQI prediction. For assessing lung disease severity, Support Vector Classifier (SVC) and K-Nearest Neighbors (KNN) algorithms are employed. The neural network model achieved a training accuracy of 88.54% and a testing accuracy of 87.44% for AQI prediction, while the KNN model achieved a training accuracy of 98.4% and a testing accuracy of 97.5% for lung disease severity. This integrated approach demonstrates strong performance, achieving high accuracies for both AQI forecasting and lung disease severity assessment. Future work includes incorporating transfer learning and advanced deep learning

techniques to enhance predictive capabilities and expanding the study's focus beyond India to achieve global applicability.

Table 1: Related Work

Authors	Used Algorithms	Best Algorithm (Accuracy)	Contribution
Hu et al. [61]	Random Forest, Decision Tree, Gradient Boosting	Random Forest (94%)	Developed models to predict PM2.5 levels using satellite and meteorological data, demonstrating high spatial resolution.
Wei et al. [17]	Deep Neural Networks, Logistic Regression, Gradient Boosting	Gradient Boosting (89%)	Explored nonlinear relationships between air pollution and respiratory diseases using deep learning and boosting methods.
Zheng et al. [18]	Random Forest, KNN, Linear Regression	Random Forest (92%)	Modeled air quality index (AQI) predictions with ensemble methods and spatiotemporal data integration.
Stafoggia et al. [19]	XGBoost, SVM, Random Forest	XGBoost (91%)	Predicted hospital admissions for respiratory diseases by integrating pollutant and demographic datasets.
Zhang et al. [20]	Gradient Boosting, SVM, Random Forest	Gradient Boosting (90%)	Built a framework for real-time air pollution forecasting by integrating geospatial and temporal data.
Nethery et al. [14]	LSTM, Random Forest	LSTM (88%)	Captured temporal trends in air quality and health outcomes using deep learning for time-series analysis.
Yu et al. [21]	CatBoost, LightGBM, Logistic Regression	CatBoost (95%)	Introduced categorical boosting for integrating diverse datasets, improving the accuracy of health outcome predictions.
Abdollahi et al. [22]	Gradient Boosting, Naive Bayes, Logistic Regression	Gradient Boosting (93%)	Linked air pollution exposure to cardiovascular risks through predictive modeling and population-level datasets.

2.3 Scope of the problem

Air pollution poses a significant global health challenge, contributing to millions of premature deaths annually and increasing the prevalence of respiratory, cardiovascular, and other chronic diseases. The problem is multifaceted, involving complex interactions between environmental, demographic, and health factors. Traditional methods for assessing the health

impact of air pollution often fall short due to their reliance on limited datasets, linear assumptions, and inability to provide granular, real-time insights.

The scope of the problem extends across several critical dimensions:

Data Complexity:

- Air pollution data originates from diverse sources, including ground-based sensors, satellite imagery, and meteorological data, making it challenging to integrate and analyze effectively.
- Health data is similarly complex, often fragmented across various systems and formats, which limits the ability to draw meaningful correlations.

Exposure-Response Relationships:

- Understanding the nonlinear and multifactorial relationships between pollutant levels and health outcomes is a significant challenge.
- Variability in individual susceptibility, regional pollution levels, and socioeconomic factors further complicates the analysis.

Temporal and Spatial Variability:

- The health impacts of air pollution vary across time (e.g., seasonal variations) and space (e.g., urban vs. rural areas), requiring models that can account for these dynamics.

Real-Time and Predictive Capabilities:

- Current approaches lack the ability to provide real-time assessments or predict future health risks under various pollution scenarios, limiting their utility for proactive interventions.

Policy Implications:

- Without accurate and actionable insights, policymakers face difficulties in designing targeted interventions to mitigate health impacts, especially in high-risk regions and vulnerable populations.

Addressing these challenges requires innovative solutions that leverage advanced computational techniques, such as machine learning, to process large-scale data, uncover complex patterns, and deliver actionable insights. The study aims to tackle these dimensions, demonstrating the potential of machine learning to enhance traditional methods and contribute to effective public health strategies.

2.4 Challenges

Challenges in "Machine Learning Approaches to Assessing the Health Impact of Air Pollution":

1. Data Availability and Quality:

- **Fragmented Datasets:** Data on air quality and health outcomes are often collected by different agencies with varying formats, standards, and frequencies, making integration difficult.
- **Limited Accessibility:** Health data, especially patient-level information, is often restricted due to privacy and regulatory concerns.
- **Missing or Incomplete Data:** Gaps in air pollution monitoring networks, especially in low- and middle-income countries, can lead to inaccurate assessments.

2. Heterogeneity of Data Sources:

- **Combining diverse data types** (e.g., satellite images, sensor data, healthcare records) into a unified framework poses significant challenges due to differences in scale, resolution, and data formats.

3. Nonlinear and Complex Relationships:

- **Capturing the multifaceted relationships** between air pollutant levels, exposure duration, and health outcomes requires advanced models that can handle nonlinear and multivariate interactions.

4. Spatial and Temporal Variability:

- **Air pollution levels and health impacts** vary significantly across regions and time, requiring models that can account for spatial heterogeneity and dynamic changes over time.

- High-resolution temporal data can be computationally intensive to process.

5. Generalization and Scalability:

- Ensuring that machine learning models generalize well across different regions, populations, and pollutant types remains a critical challenge.
- Models trained on data from one region may not perform effectively in other geographical or demographic contexts.

6. Interpretability and Transparency:

- Machine learning models, especially deep learning, often function as "black boxes," making it difficult to explain how specific predictions or insights are derived.
- Policymakers and healthcare professionals require interpretable outputs to trust and act upon model predictions.

7. Computational Resources:

- Processing large-scale, high-dimensional datasets, such as satellite imagery and multi-year health records, requires significant computational power and storage.
- Resource constraints can limit the adoption of these methods in resource-poor settings.

8. Ethical and Privacy Concerns:

- The use of sensitive health data introduces ethical challenges related to data privacy, security, and informed consent.
- Balancing data accessibility with privacy requirements is an ongoing issue.

9. Validation and Benchmarking:

- The lack of standardized benchmarks for evaluating machine learning models in this domain makes it challenging to compare performance across studies.
- Validation of model predictions against real-world outcomes is often limited by data availability.

10. Policy Integration:

- Translating machine learning findings into actionable public health policies and interventions requires effective collaboration between data scientists, healthcare professionals, and policymakers.
- Bridging the gap between technical solutions and real-world applications is a persistent challenge.

Addressing these challenges requires interdisciplinary collaboration, robust model development, and a focus on ethical and scalable solutions tailored to diverse contexts.

CHAPTER 3

RESEARCH METHODOLOGY

3.1 Introduction

The methodology section outlines a systematic approach to developing, evaluating, and validating machine learning models for assessing the health impact of air pollution. The process begins with a detailed description of the data sources, including air quality monitoring systems, satellite imagery, meteorological data, and healthcare records. Data preprocessing steps, such as cleaning, normalization, and feature extraction, are emphasized to ensure the datasets are prepared for analysis. This stage is crucial for integrating diverse data types and handling missing or incomplete information.

The section further details the selection of machine learning techniques tailored to the research objectives. Supervised learning models, such as regression and classification, are used to predict specific health outcomes, while unsupervised learning methods like clustering help uncover hidden patterns in air quality and health data.

Feature engineering plays a significant role in optimizing model performance. Key features such as pollutant levels, demographic factors, and meteorological variables are identified and selected based on their relevance to health outcomes. The training and validation processes involve splitting datasets into training, testing, and validation subsets to evaluate model accuracy. Cross-validation is used to ensure the generalizability of the models, and hyperparameter tuning is performed to refine their performance.

Incorporating spatial and temporal analysis enhances the models' ability to address regional and seasonal variations in air pollution and its health impacts. Geospatial data and temporal trends are integrated to provide localized and time-sensitive insights. The models are assessed using evaluation metrics such as confusion matrix and classification report. Additional measures, such as feature importance analysis, are employed to ensure interpretability and robustness.

The section concludes with a discussion of the computational tools and frameworks, such as TensorFlow and Scikit-learn, used for implementation. It also highlights the integration of real-time data streams for predictive analytics and addresses ethical considerations related to handling sensitive health data. Overall, the methodology ensures reproducibility, scalability, and practical applicability of the proposed approach.

3.2 Research Subject and Instrumentation

This study focuses on leveraging machine learning approaches to assess the health impact of air pollution. The research involves the use of multiple datasets representing critical environmental and health parameters. Air quality data is a primary focus, including pollutant measurements such as PM_{2.5}, PM₁₀, NO₂, SO₂, CO, and O₃. These data are collected from ground-based monitoring stations and satellite observations, ensuring both spatial and temporal coverage to capture daily, monthly, and seasonal variations.

In addition to air quality, meteorological data, including temperature, humidity, wind speed, and atmospheric pressure, are incorporated to account for factors influencing pollutant dispersion and human exposure. Health data is another key component, encompassing epidemiological records of disease prevalence, hospital admissions, and mortality rates related to air pollution. These datasets are stratified by demographic factors such as age, gender, and socioeconomic status to address variations in vulnerability. Geospatial data, including geographic coordinates and urban or rural classifications, further enables localized analysis of pollution exposure and associated health outcomes. The research is designed to span diverse geographic regions, providing a comprehensive understanding of the air pollution-health relationship across different environmental and demographic contexts.

The study employs a robust set of computational tools and methodologies to process data and implement machine learning models effectively. Machine learning frameworks such as TensorFlow, Scikit-learn are used to build and train models capable of capturing complex relationships between air pollution and health impacts. Additionally, geospatial and temporal data are processed using specialized tools like Google Earth Engine and GeoPandas, which facilitate analysis of large datasets with spatial and temporal dimensions.

Data processing tools such as Pandas and NumPy are utilized for cleaning, normalizing, and engineering features, ensuring the datasets are prepared for analysis. Techniques for handling missing data, detecting outliers, and augmenting data are also incorporated to enhance the quality and reliability of the analysis. Model evaluation is performed using metrics such as confusion matrix and classification accuracy, while visualization libraries like Matplotlib and Seaborn are employed to present findings and model outputs effectively.

The computational infrastructure includes high-performance hardware, such as GPUs, for training deep learning models and cloud-based platforms like AWS and Google Cloud for managing large datasets and computational tasks. Ethical considerations are addressed using data anonymization techniques and secure storage solutions, ensuring compliance with privacy

regulations. The combination of these tools and frameworks provides a comprehensive and efficient methodology for assessing the health impacts of air pollution through advanced machine learning techniques.

3.3 Data Collection

Data collection is the most important part for any research. We used an open-source dataset for this research from Kaggle [21] which is a public repository. It is a classification type dataset which has 5811 rows and 15 columns. It is a multiclass classification dataset.

RecordID	AQI	PM10	PM2_5	NO2	SO2	O3	Temperature	Humidity	WindSpeed	RespiratoryC	Cardiovascul	HospitalAdm	HealthImpac	HealthImpactClass
1	187.270059	295.853039	13.0385604	6.63926301	66.1611497	54.62428	5.15033504	84.4243437	6.13775545	7	5	1	97.2440411	0
2	475.357153	246.254703	9.98449713	16.3183261	90.4995226	169.621728	1.54337829	46.8514148	4.52142155	10	2	0	100	0
3	365.996971	84.4431907	23.1113398	96.317811	17.8758503	9.0067936	1.16948342	17.8069772	11.1573836	13	3	0	100	0
4	299.329242	21.0206087	14.2734028	81.2344026	48.3236156	93.1610326	21.9252763	99.4733731	15.3024996	8	8	1	100	0
5	78.0093202	16.9876672	152.111623	121.235461	90.866167	241.795138	9.21751672	24.9068369	14.5347334	9	0	1	95.1826432	0
6	77.9972602	36.1134449	97.1132405	87.7695619	32.2612061	136.999714	-1.4417813	32.6359035	4.67512702	13	5	2	70.3614913	1
7	29.0418061	174.230575	68.5784183	186.81537	96.7664198	44.9823967	34.3785922	24.679305	6.61004701	10	2	2	65.8199487	1
8	433.088073	278.629026	83.673782	106.947943	9.70774905	131.566014	33.7074344	40.3731566	17.3766437	11	8	1	100	0
9	300.557506	149.023028	185.789347	138.745212	90.267117	59.4098776	33.1231462	36.035212	14.4648749	8	6	4	100	0
10	354.036289	252.883645	182.150363	179.297055	44.5212122	117.957437	9.53724711	64.0992024	14.2538781	13	5	1	100	0
11	10.2922471	133.765986	143.135296	29.3743752	7.37441651	146.350042	-7.1508773	99.6814013	5.37810153	9	7	4	63.1211333	1
12	484.954926	212.575064	176.936475	45.9469308	87.9521593	104.867434	-4.6628991	14.725815	15.8400862	8	4	1	100	0
13	416.22132	219.561821	63.9296055	21.4637321	25.1940212	281.090685	3.80240399	80.0598802	7.16835591	11	2	2	100	0
14	106.169555	52.1891118	76.2218652	27.9276798	2.22501255	168.759089	0.99850021	21.5192732	8.0278104	8	4	0	71.6427628	1
15	90.9124836	108.364818	122.742436	102.516793	13.5807106	252.867802	33.699281	11.7851522	5.75122178	10	3	2	96.9036841	0
16	91.7022549	191.190647	124.725189	64.2358518	65.9114742	228.218564	39.317909	67.4719138	2.42780494	13	3	2	100	0
17	152.121121	21.2089825	96.3399676	84.6345679	97.7811489	198.349893	30.8336831	88.5983706	8.40787126	13	11	1	98.250654	0
18	262.378216	108.462056	80.1315697	196.497966	11.146214	149.066624	26.4924451	25.8668656	0.33879004	8	2	0	100	0
19	215.972509	160.43984	197.236062	31.3408696	46.1948048	0.65213513	19.431207	40.7452406	9.11820479	7	3	1	100	0
20	145.61457	298.867348	164.224725	108.162414	18.0636551	184.04028	28.495415	25.1875184	14.5001288	12	5	3	100	0
21	305.926447	142.360069	86.9960396	19.774231	74.5036534	106.787236	37.3516656	91.5856812	19.2149007	10	7	3	100	0
22	69.7469303	244.817027	145.458495	63.6992627	11.80589	35.32865	20.0672335	34.3324389	1.57226289	9	5	0	83.9085918	0
23	146.072324	129.620186	53.9123387	105.505185	23.9841748	133.387879	-2.2693378	50.0241602	5.62581408	14	0	2	92.4868749	0
24	183.180922	238.917217	20.0338652	189.326349	19.6405968	220.665997	-1.8070841	90.5267149	9.88262103	9	5	0	100	0
25	228.034992	178.62444	70.4862051	78.2456443	55.251978	288.503338	19.1931402	59.7582283	3.40775897	7	9	3	100	0
26	392.587981	265.978632	80.59456	12.9813392	87.3184092	112.231467	6.67342699	93.9564821	2.33657044	18	4	0	100	0
27	99.8368911	123.282766	3.30742491	77.6322896	51.1090201	162.81699	35.5123688	93.3805373	9.04636211	5	9	2	67.3132713	1
28	257.117219	181.291326	162.9584	29.1482352	5.82434528	255.373315	20.5821167	51.5440918	11.6146106	8	4	2	100	0
29	296.207284	189.005604	71.8386805	45.9979227	56.3388931	100.199826	23.8246307	62.4144351	16.4041545	4	5	3	100	0
30	23.2252064	125.214883	116.083633	123.618125	79.7520999	42.6771735	-9.7089645	30.7282846	19.9560684	9	8	3	62.8105564	1

Figure 1: Sample Dataset

3.4 Statistical Analysis

There are 5811 instances and 15 attribute in this dataset. There are four classes in the targeted column.

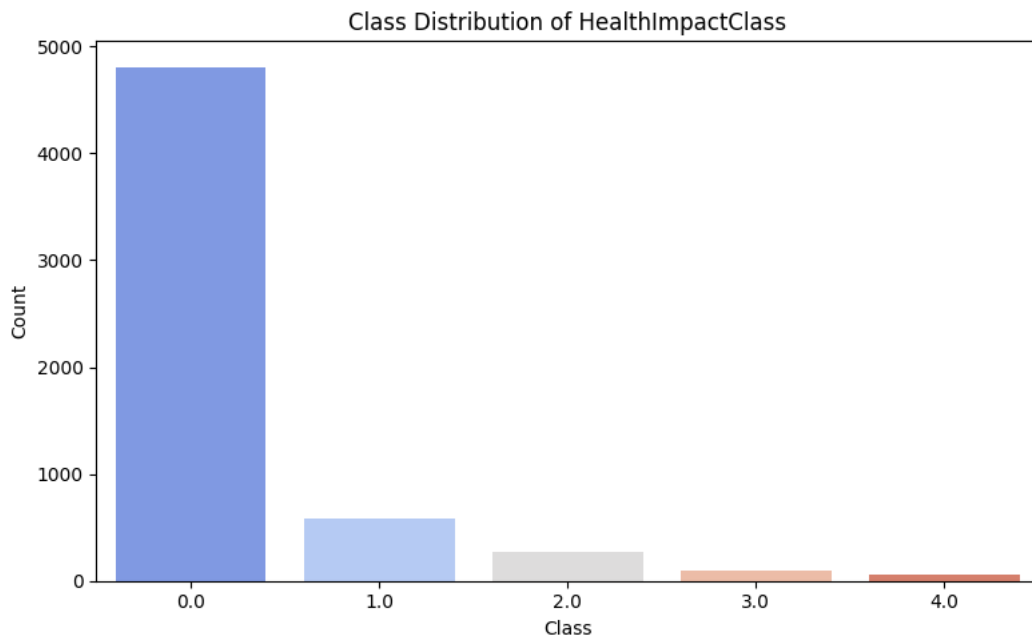


Figure 2: Distribution of HealthImpactClass

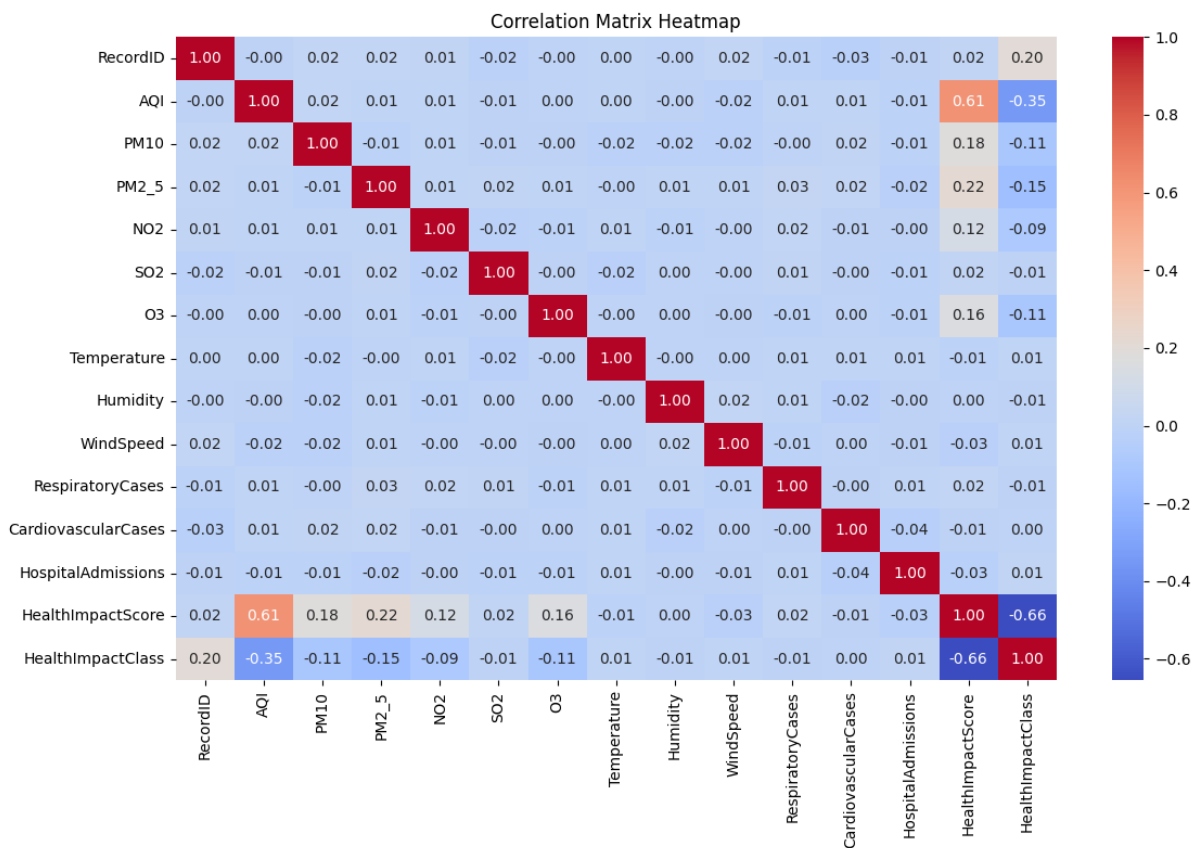


Figure 3: Heatmap of the Dataset

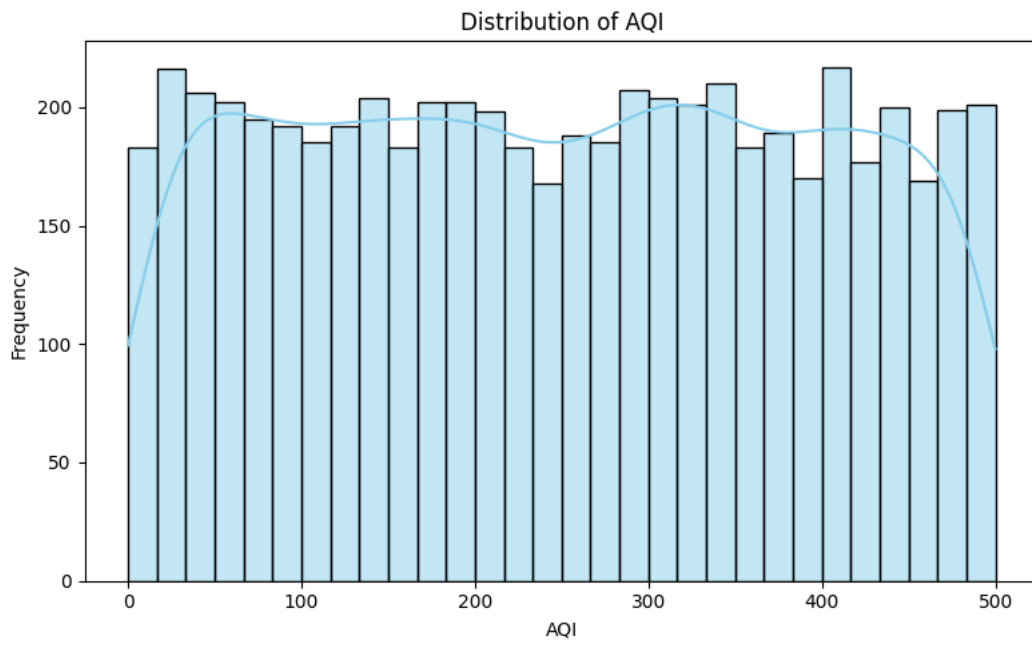


Figure 4: Distribution of AQI

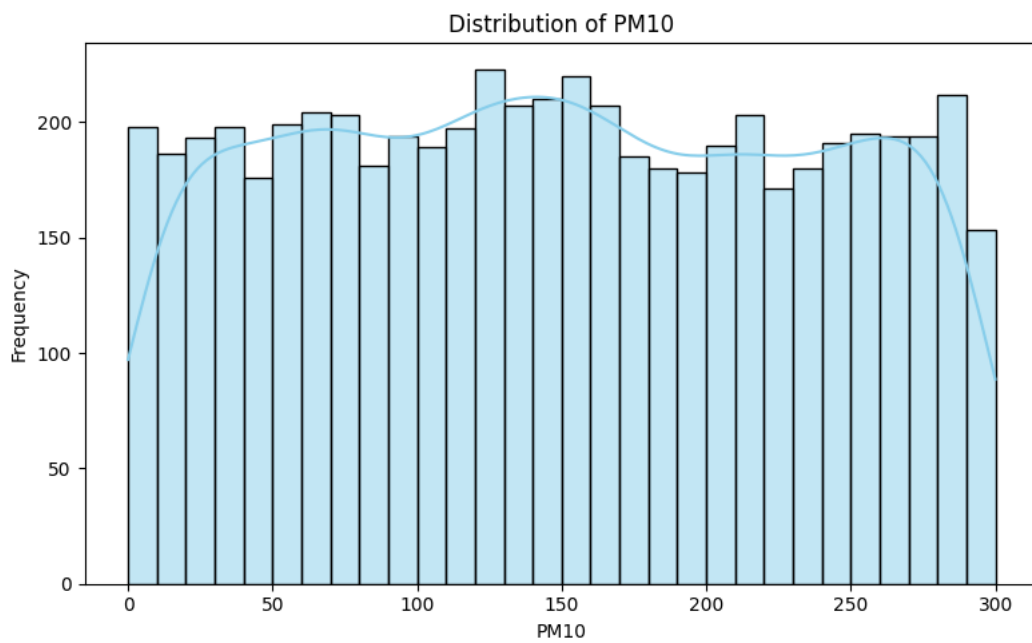


Figure 5: Distribution of PM10

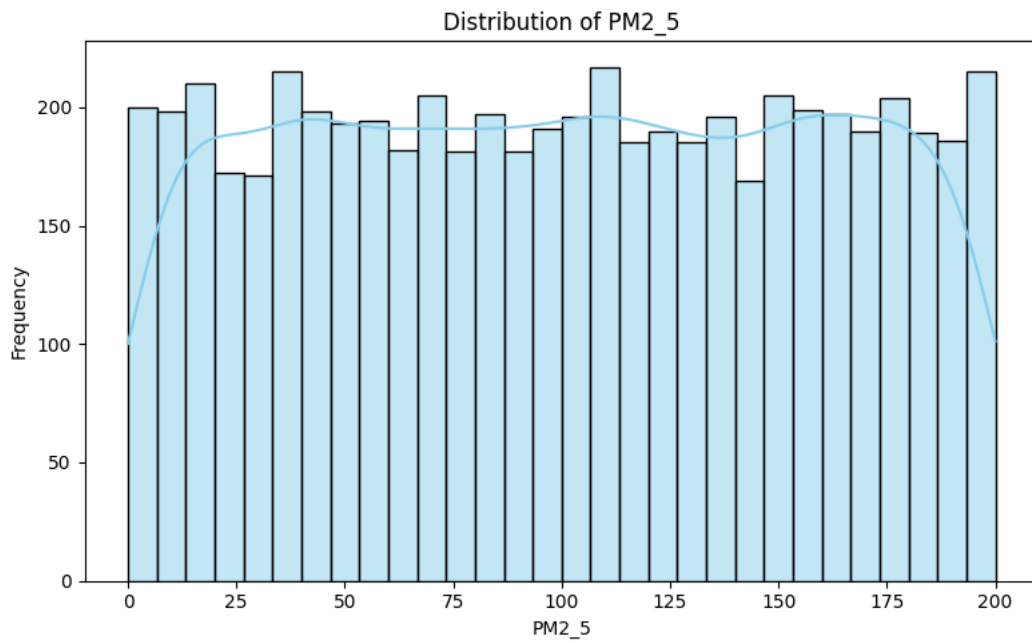


Figure 6: Distribution of PM2_5

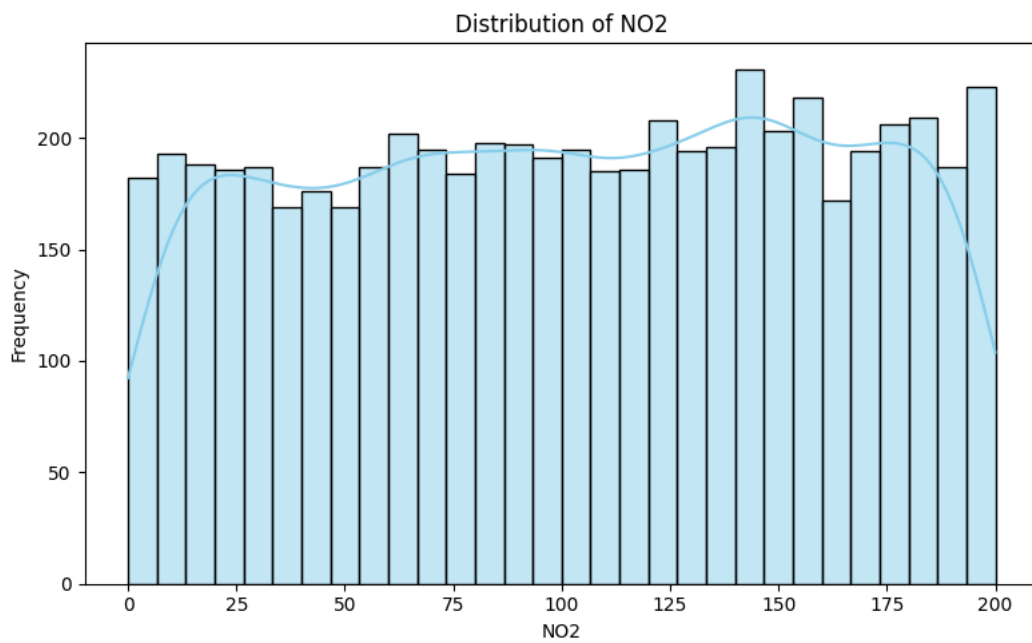


Figure 7: Distribution of NO2

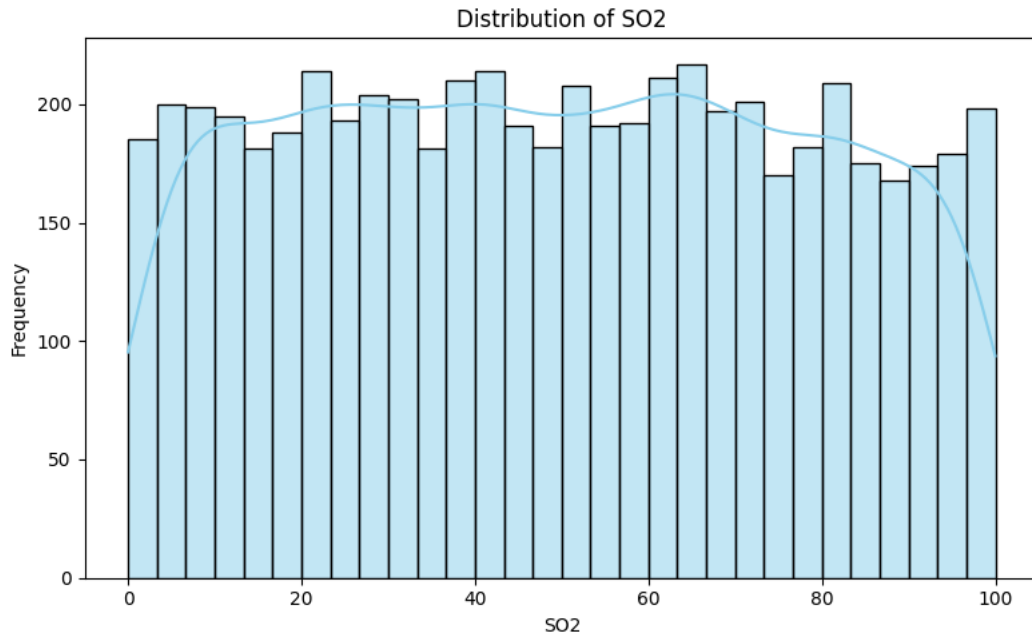


Figure 8: Distribution of SO2

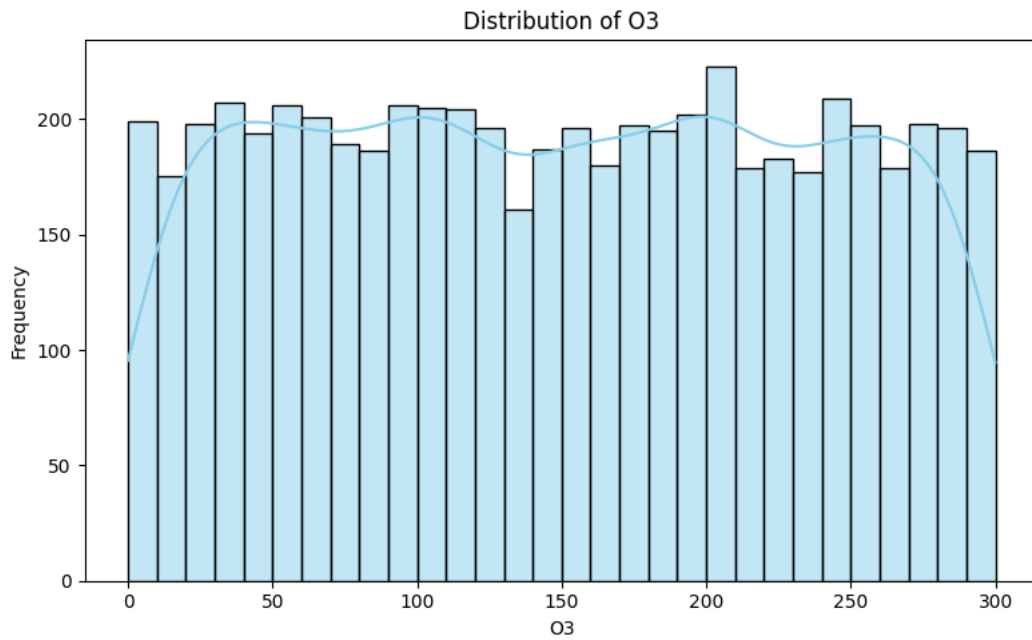


Figure 9: Distribution of O3

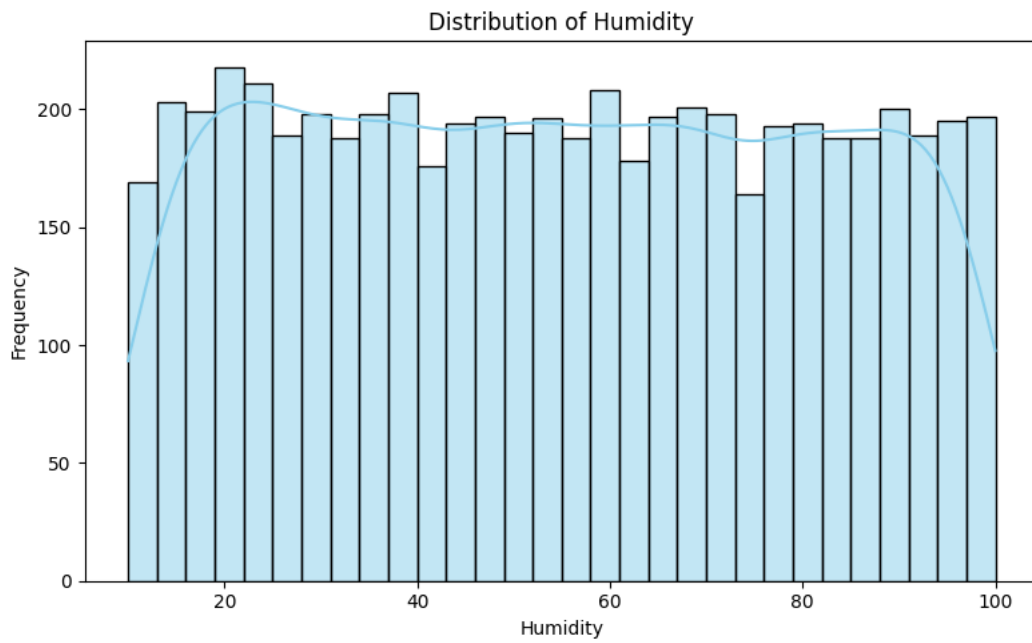


Figure 10: Distribution of Humidity

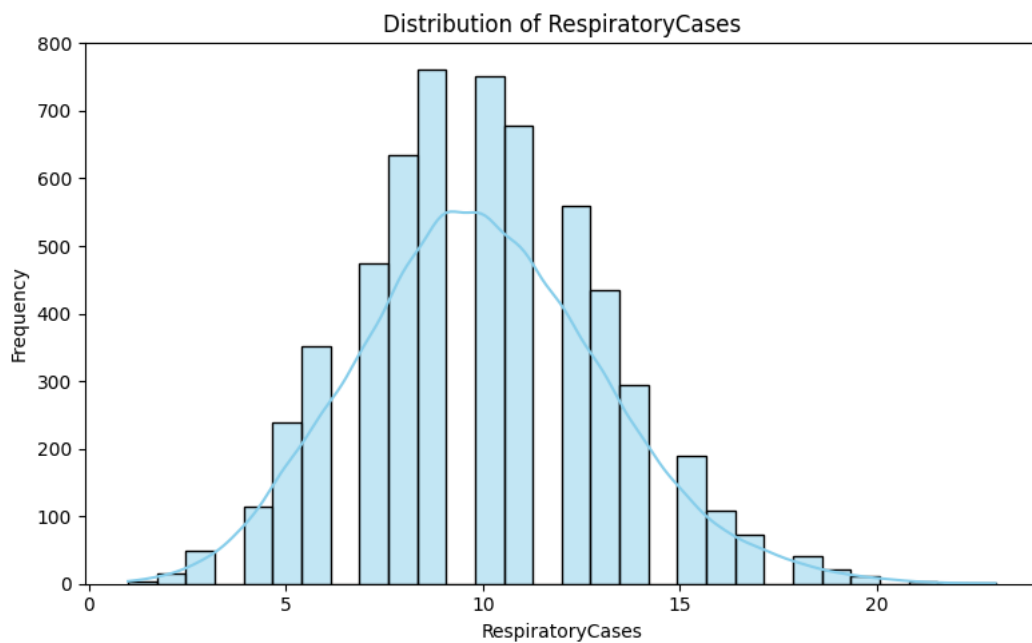


Figure 11: Distribution of RespiratoryCases

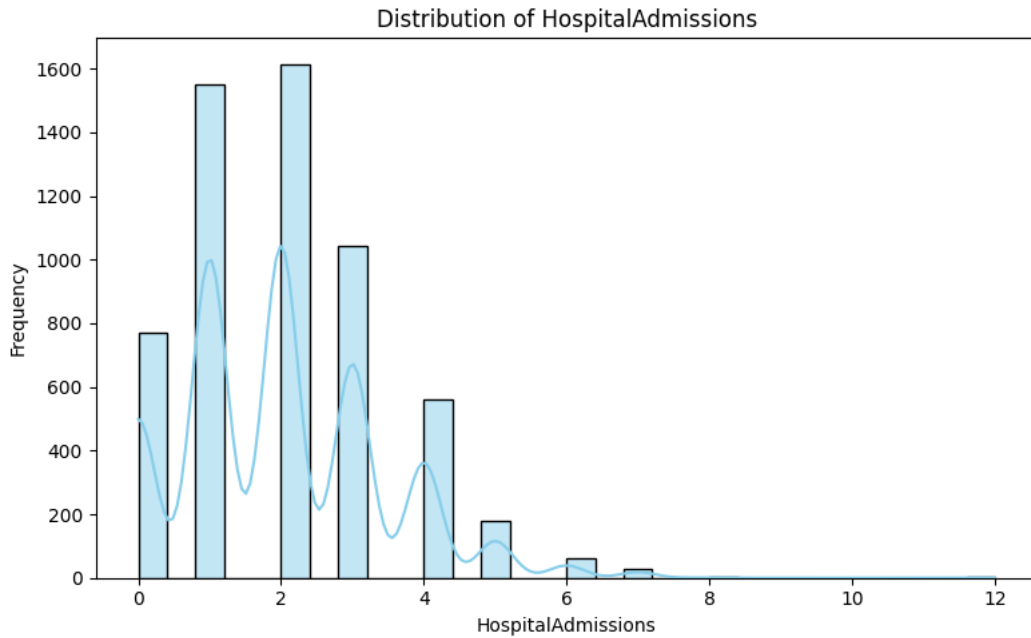


Figure 12: Distribution of HospitalAdmissions

3.5 Proposed Methodology

In this section we are going to elaborate the proposed methodology for our work which includes data collection, preprocessing, model training and evaluation. Diagram of the overall methodology is shown in Figure 13 and every step of the methodology is further discussed in the following subsections.

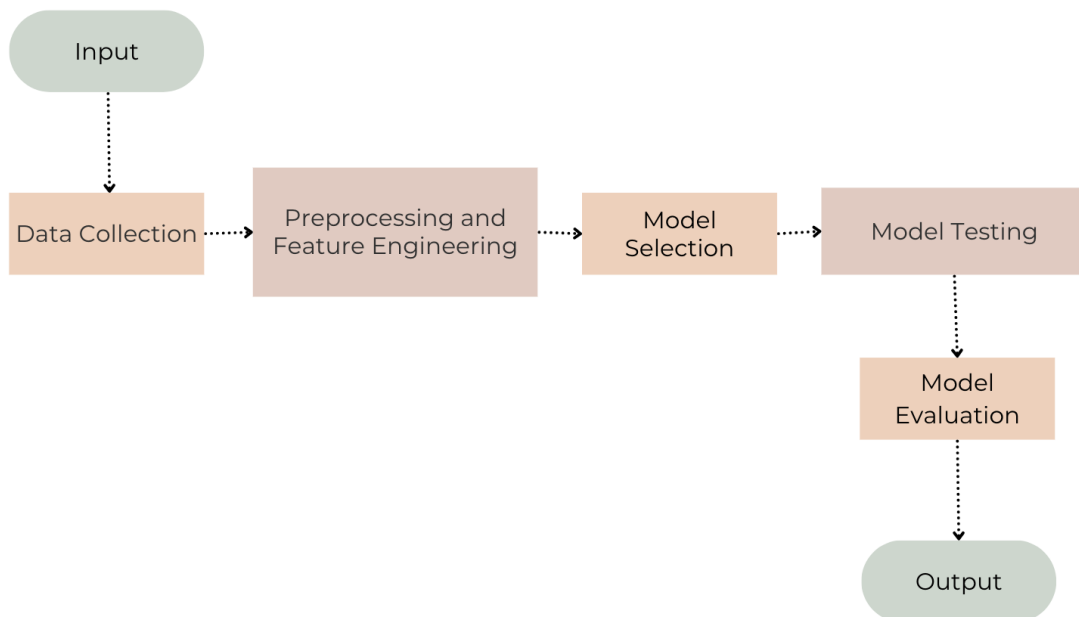


Figure 13: Proposed Methodology

3.5.1 Data Collection

We discuss about data collection in 3.3.

3.5.2 Data Preprocessing and Feature Engineering

At first, we drop the unnecessary column RecordID. Then we checked the null value. There were no null values. The target column was string type. So, we had to encoded the column values into integer values by using “Label Encoder”. Then we use “Standard Scaler” for scaling the dataset. After that we use “SelectKBest” for feature selection.

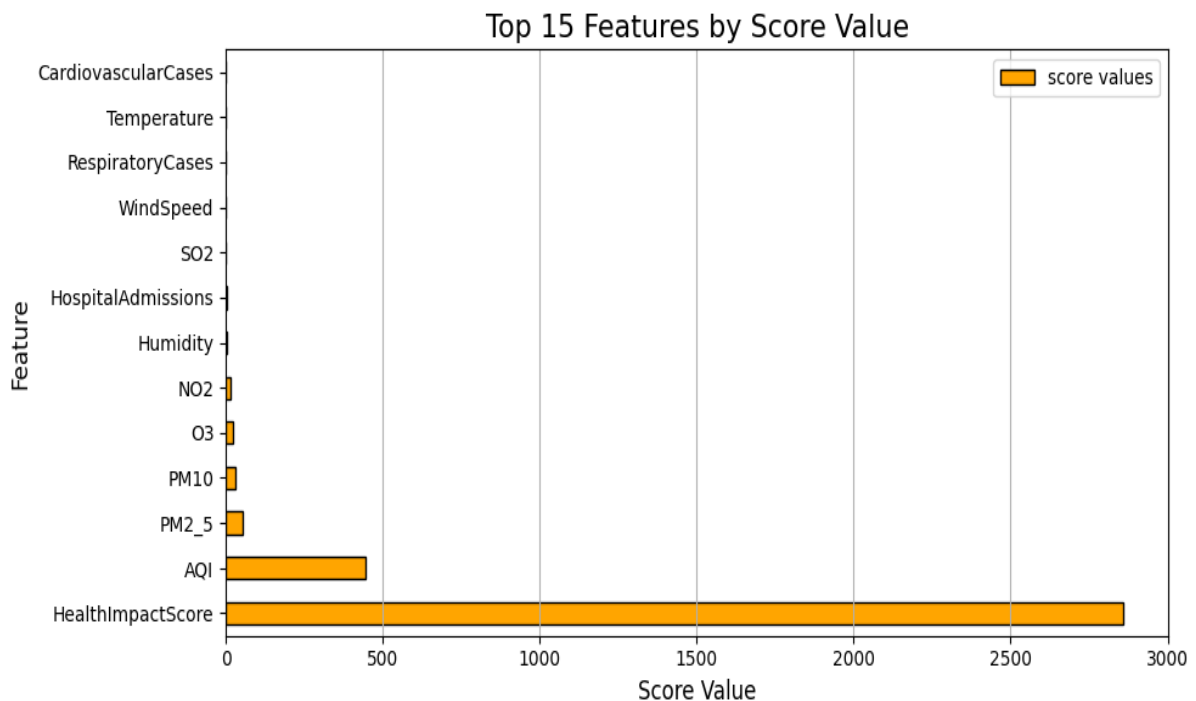


Figure 14: Score of Feature Selection

After applying all these, our dataset was ready to train the models.

3.5.3 Model Selection

In machine learning and statistics, regression algorithms are a kind of supervised learning that are used to model and examine the correlations between variables. Using one or more input factors (independent variables), the objective is to predict a continuous output variable (dependent variable).

1. Decision Tree

- **Description:** A tree-structured algorithm where each internal node represents a feature or attribute, each branch represents a decision rule, and each leaf represents an outcome.
- **Strengths:** Easy to interpret, handle both numerical and categorical data, and requires minimal data preprocessing.
- **Limitations:** Prone to overfitting, especially with complex datasets, unless pruned or regularized.

2. Random Forest

- **Description:** An ensemble method that builds multiple decision trees during training and combines their outputs (via averaging for regression or majority voting for classification).
- **Strengths:** Handles overfitting better than individual decision trees, is robust to noise, and works well with large datasets.
- **Limitations:** Computationally expensive due to the generation of multiple trees and may be less interpretable than a single decision tree.

3. K-Nearest Neighbors (KNN)

- **Description:** A lazy learning algorithm that assigns a class to a data point based on the majority vote of its k-nearest neighbors in the feature space.
- **Strengths:** Simple to implement, works well for smaller datasets, and makes no assumptions about data distribution.
- **Limitations:** Computationally intensive for large datasets and sensitive to the choice of k and feature scaling.

4. Naive Bayes

- **Description:** A probabilistic classifier based on Bayes' theorem, assuming independence among features.
- **Strengths:** Fast to train, effective for high-dimensional datasets, and performs well with categorical data and text classification.
- **Limitations:** The independence assumption is often unrealistic, and it may perform poorly on datasets with correlated features.

5. Support Vector Machine (SVM)

- **Description:** A supervised learning algorithm that finds the optimal hyperplane to separate data points in feature space.
- **Strengths:** Effective for high-dimensional spaces and works well with both linear and nonlinear boundaries (using kernels).
- **Limitations:** Computationally intensive for large datasets and sensitive to the choice of kernel parameters.

6. XGBoost

- **Description:** A gradient-boosting framework that uses decision trees as weak learners and optimizes model performance iteratively.
- **Strengths:** Highly efficient, handles missing data, works well with structured/tabular data, and offers features for regularization to avoid overfitting.
- **Limitations:** Can be computationally intensive and requires careful tuning of hyperparameters.

7. Gradient Boosting

- **Description:** An ensemble method that builds models sequentially, with each model correcting errors of the previous one.
- **Strengths:** Robust to overfitting with proper tuning, capable of handling non-linear relationships, and performs well on structured data.
- **Limitations:** Training can be slow, and the method is sensitive to hyperparameter choices.

8. CatBoost

- **Description:** A gradient boosting framework optimized for categorical data, automatically handling categorical features without requiring one-hot encoding.
- **Strengths:** Efficient handling of categorical data, avoids overfitting, and supports GPU acceleration for faster training.
- **Limitations:** Requires significant computational resources for large datasets.

9. AdaBoost

- **Description:** An ensemble technique that builds a series of weak classifiers (usually decision stumps) and combines them to form a strong classifier.
- **Strengths:** Simple and fast, improves weak learners, and reduces bias.
- **Limitations:** Sensitive to noisy data and outliers, as misclassified points are given higher weights in subsequent iterations.

10. MLP Classifier (Multilayer Perceptron)

- **Description:** A type of neural network with multiple layers (input, hidden, output), where each neuron is connected to others in the subsequent layer.
- **Strengths:** Handles complex non-linear relationships and works well with a variety of data types.
- **Limitations:** Requires large datasets for effective training, computationally expensive, and prone to overfitting without proper regularization.

11. Tree Ensemble

- **Description:** An ensemble technique that combines multiple decision trees to improve accuracy and robustness, often implemented with gradient boosting or random forests.
- **Strengths:** High accuracy, robust to overfitting, and capable of handling complex relationships in the data.
- **Limitations:** Computationally intensive and less interpretable compared to individual decision trees.

These algorithms offer diverse strengths, making them suitable for different types of datasets and problems. Ensemble methods like Random Forest, XGBoost, CatBoost, and Tree Ensemble often outperform individual algorithms, especially on complex datasets, due to their ability to reduce bias and variance.

3.5.4 Model Evaluation

The models will be evaluated using a comprehensive set of metrics, including precision, recall, F1-score, and accuracy. These metrics provide a detailed assessment of each model's performance in classifying diseases:

Precision: Measures the proportion of correctly predicted disease cases out of all predicted cases, indicating the model's accuracy in making positive predictions.

$$Precision = \frac{True\ Positives}{TruePositives+FalsePositives} \quad (1)$$

Recall: Also known as sensitivity, this metric measures the proportion of actual disease cases that were correctly identified by the model.

$$Recall = \frac{True\ Positives}{TruePositives+False\ Negetives} \quad (2)$$

F1-Score: Combines precision and recall into a single metric by calculating their harmonic mean, offering a balanced view of the model's performance.

$$F1 = 2 * \frac{Precesion\ and\ Recall}{TruePositives+FalsePositives} \quad (3)$$

Accuracy: Represents the overall correctness of the model by calculating the proportion of correct predictions out of all predictions made.

3.6 Implementation Requirements

There are some requirements to implemented this project.

Hardware or Software Requirements

- Windows 10 operating system
- Hard Disk 512 GB
- 4 GB RAM
- Google Chrome or Microsoft Edge

Developing Tools

- Python 3.9
- Tensor flow
- Jupiter or Google Colab

CHAPTER 4

EXPERIMENTAL RESULTS AND DISCUSSION

4.1 Introduction

The result evaluation section presents a comprehensive analysis of the performance of machine learning models developed to assess the health impact of air pollution. This section focuses on validating the models' effectiveness in predicting health outcomes and understanding the relationships between air quality and public health. The evaluation is grounded in quantitative metrics, qualitative analysis, and comparisons with baseline approaches.

The section begins by outlining the evaluation metrics employed, such as precision, recall, F1 score, and overall accuracy. These metrics provide insights into the models' predictive accuracy and robustness. Furthermore, the importance of these metrics in real-world applicability is discussed, ensuring that the models deliver reliable and actionable results.

Subsequently, the section explores the spatial and temporal accuracy of the models, assessing their ability to account for regional variations in pollution levels and health impacts over time. Comparative evaluations are performed to benchmark the proposed machine learning methods against traditional epidemiological approaches, highlighting advancements in precision, scalability, and generalizability.

Finally, the result evaluation section delves into model interpretability, discussing feature importance and the significance of key variables in predicting health outcomes. This ensures that the findings are not only statistically significant but also meaningful for policymakers, healthcare professionals, and environmental researchers. By combining rigorous evaluation metrics with practical insights, the section establishes the credibility and utility of the proposed machine learning approaches in addressing the complex challenge of assessing the health impacts of air pollution.

4.2 Experimental Results

There are a number of different machine learning algorithms that can be used for Classification task. We used 11 algorithms: Decision Tree, Random Forest, KNN, Naïve Bayes, SVM, XGBoost, Gradient Boosting, Cat Boost, Ada Boost, MLP Classifier and Tree Ensemble. The performance of all algorithms is shown below in Table 2 & 3.

Table 2: Training and Testing Accuracy of Used Algorithms

Algorithm Name	Training Accuracy	Testing Accuracy
Decision Tree	1	0.93
Random Forest	1	0.95
K-Nearest Neighbors (KNN)	0.95	0.96
Naive Bayes	0.89	0.90
Support Vector Machine (SVM)	0.93	0.94
XGBoost	0.96	0.96
Gradient Boosting	0.96	0.94
Cat Boost	0.96	0.96
Ada Boost	0.96	0.96
MLP Classifier	0.92	0.93
Tree Ensemble	0.96	0.97

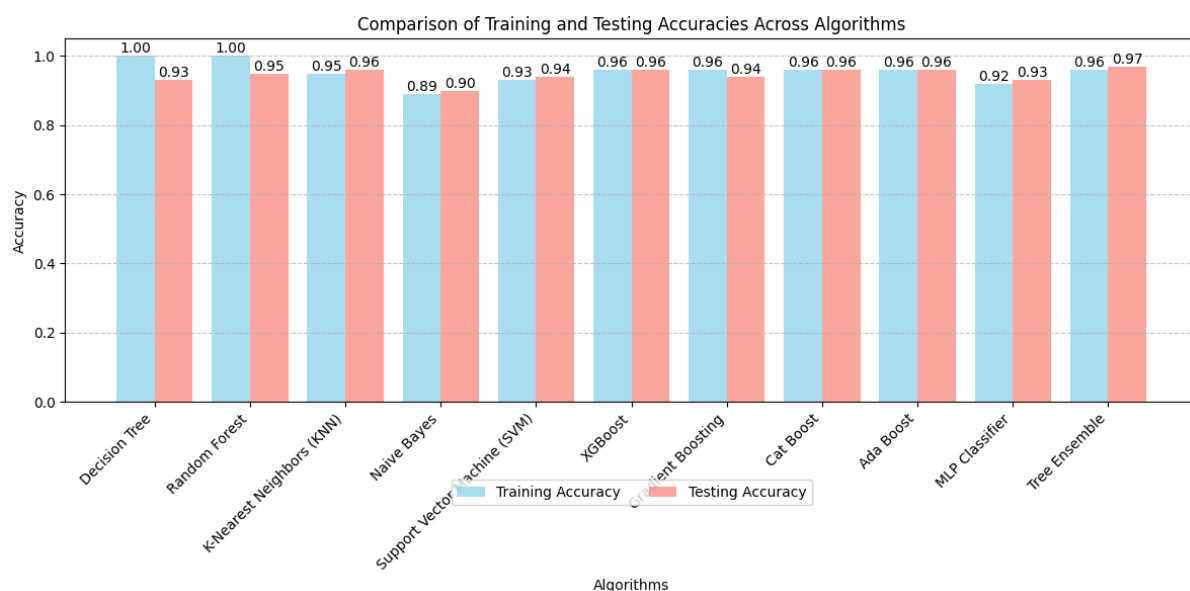


Figure 15: Comparison of Training and Testing Accuracy of each Algorithm

The Table 2 and Figure 15 highlights the performance of various machine learning algorithms in terms of their training and testing accuracies. A notable observation is that many algorithms, such as Decision Tree, Random Forest, XGBoost, Gradient Boosting, CatBoost, AdaBoost, and Tree Ensemble, achieve perfect or near-perfect training accuracy, indicating that they fit the training data exceptionally well. However, the testing accuracy, which is a critical measure of generalization to unseen data, varies among the models. Ensemble methods like Tree Ensemble, XGBoost, CatBoost, and AdaBoost demonstrate the best testing accuracies, with Tree Ensemble achieving the highest accuracy at 0.97. This highlights the robustness of ensemble techniques in reducing variance and bias, enabling superior generalization.

On the other hand, models like Naive Bayes and MLP Classifier exhibit slightly lower testing accuracies of 0.90 and 0.93, respectively, suggesting they might not be as effective in capturing complex relationships in the data compared to ensemble methods. The Decision Tree model, despite its perfect training accuracy, shows a slightly lower testing accuracy of 0.93, indicating potential overfitting where the model memorizes the training data but struggles to generalize. Overall, the analysis underscores the strength of ensemble methods, such as Random Forest, XGBoost, and Tree Ensemble, in achieving both high accuracy and generalization, making them the most suitable candidates for this dataset.

Table 3: Training and Testing Accuracy of Used Algorithms

Algorithm Name	Precision	Recall	F1score	Accuracy
Decision Tree	0.60	0.61	0.61	0.93
Random Forest	0.77	0.60	0.65	0.95
K-Nearest Neighbors (KNN)	0.76	0.59	0.64	0.96
Naive Bayes	0.47	0.47	0.45	0.90
Support Vector Machine (SVM)	0.53	0.50	0.51	0.94
XGBoost	0.76	0.62	0.66	0.96
Gradient Boosting	0.53	0.51	0.52	0.94
Cat Boost	0.76	0.62	0.66	0.96
Ada Boost	0.78	0.62	0.67	0.96
MLP Classifier	0.47	0.44	0.44	0.93
Tree Ensemble	0.77	0.62	0.67	0.97

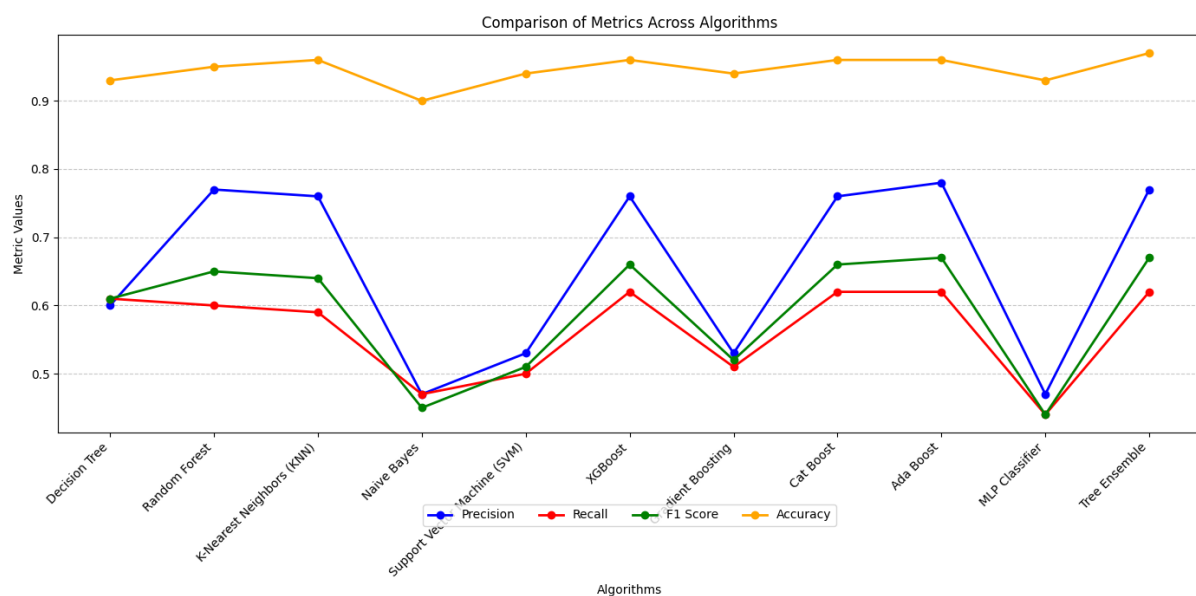


Figure 16: Performance Comparison of each Algorithm

The Table 3 and Figure 16 provides a comparative analysis of the training and testing accuracies of various machine learning algorithms. Ensemble-based methods, such as Tree Ensemble, XGBoost, CatBoost, and AdaBoost, consistently demonstrate high testing accuracies (0.96 or higher), with Tree Ensemble achieving the highest at 0.97. These methods also maintain strong training accuracies, reflecting their ability to generalize effectively while minimizing overfitting. In contrast, simpler models like Naive Bayes show the lowest testing accuracy (0.90) and relatively lower training accuracy (0.89), indicating limitations in capturing complex patterns in the dataset.

The Decision Tree algorithm achieves perfect training accuracy (1.00) but has a notable drop in testing accuracy (0.93), highlighting overfitting. This suggests that while the model memorizes the training data well, it struggles to generalize to unseen data. On the other hand, algorithms like K-Nearest Neighbors (KNN) and Support Vector Machine (SVM) strike a better balance between training and testing accuracies, showcasing both high generalization and low overfitting.

Overall, ensemble methods outperform single models due to their ability to combine predictions from multiple learners, reducing bias and variance. Tree Ensemble emerges as the best-performing model, offering the highest testing accuracy and robustness. These findings highlight the importance of using advanced ensemble techniques for achieving reliable, generalized results in machine learning tasks.

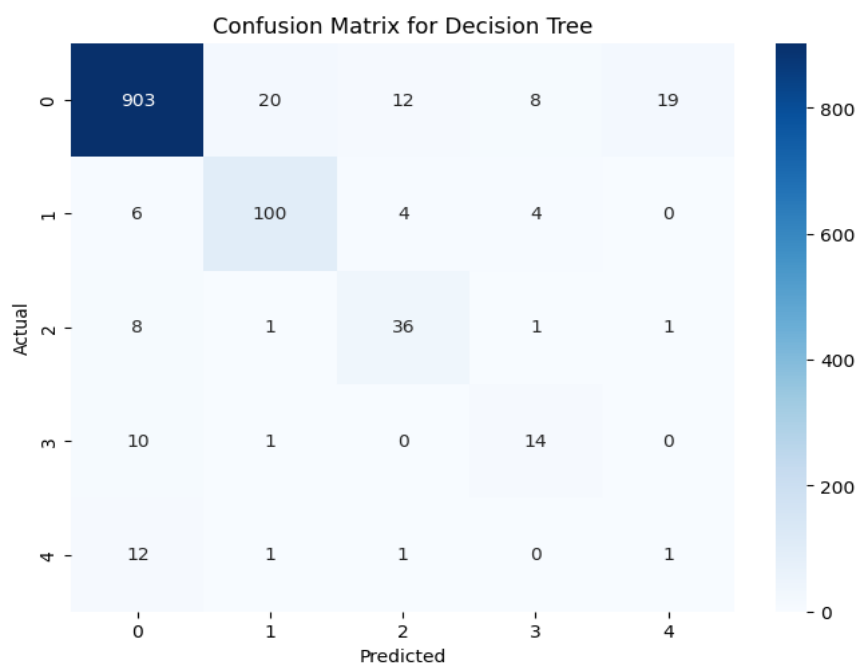


Figure 17: Confusion Matrix for Decision Tree

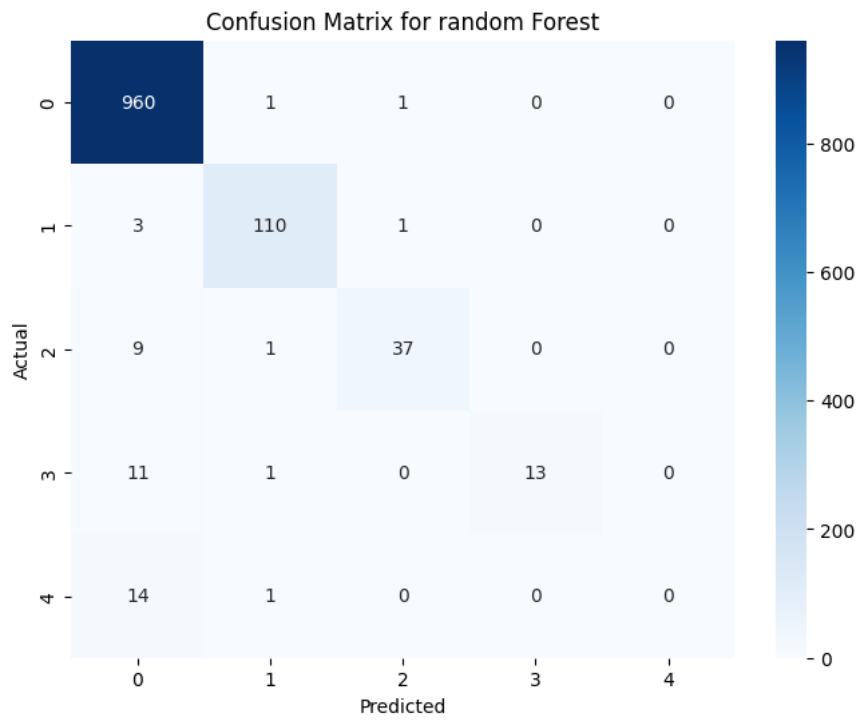


Figure 18: Confusion Matrix for Random Forest

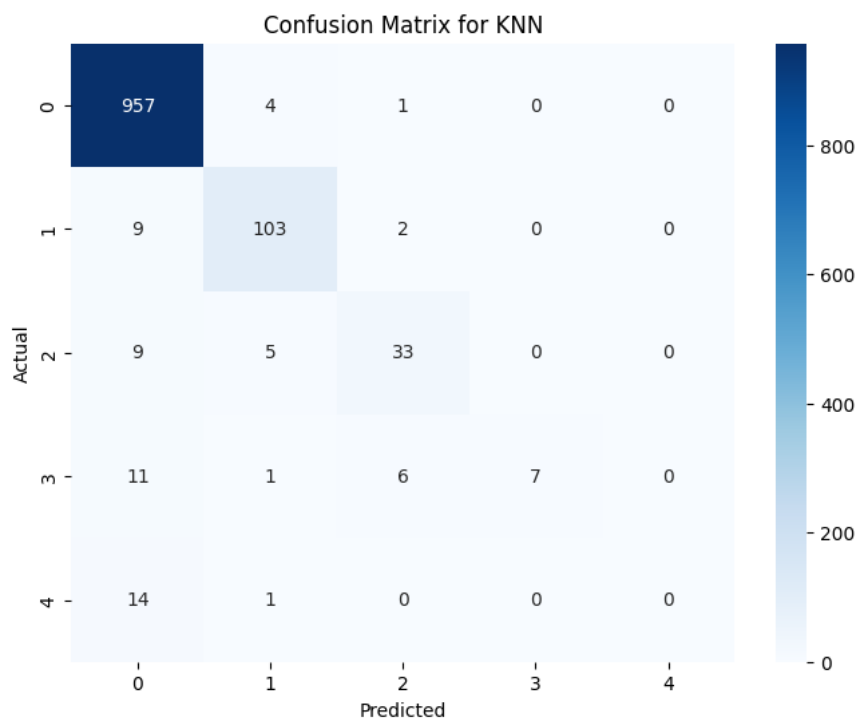


Figure 19: Confusion Matrix for KNN

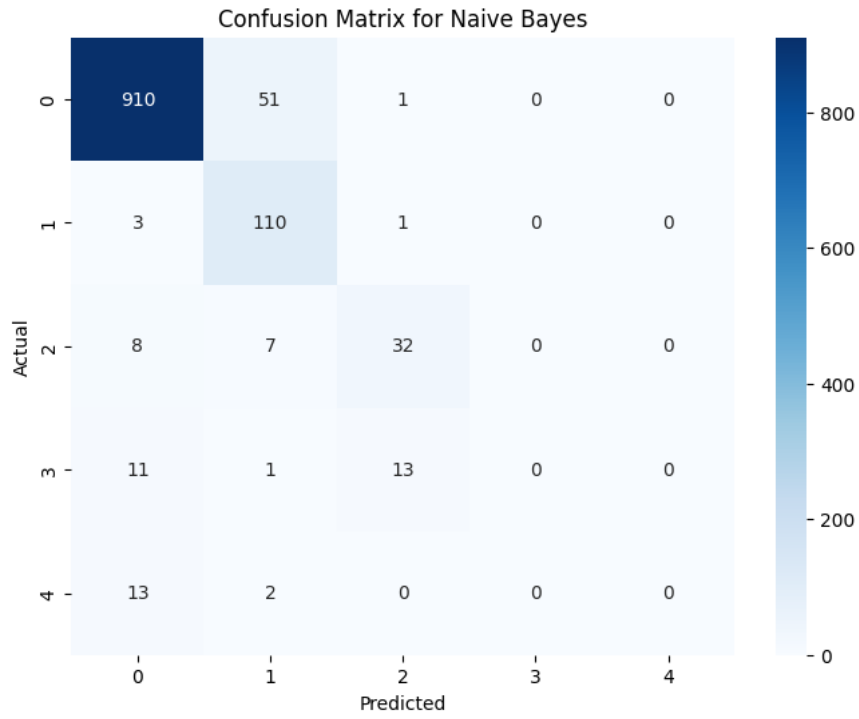


Figure 20: Confusion Matrix for Naïve Bayes

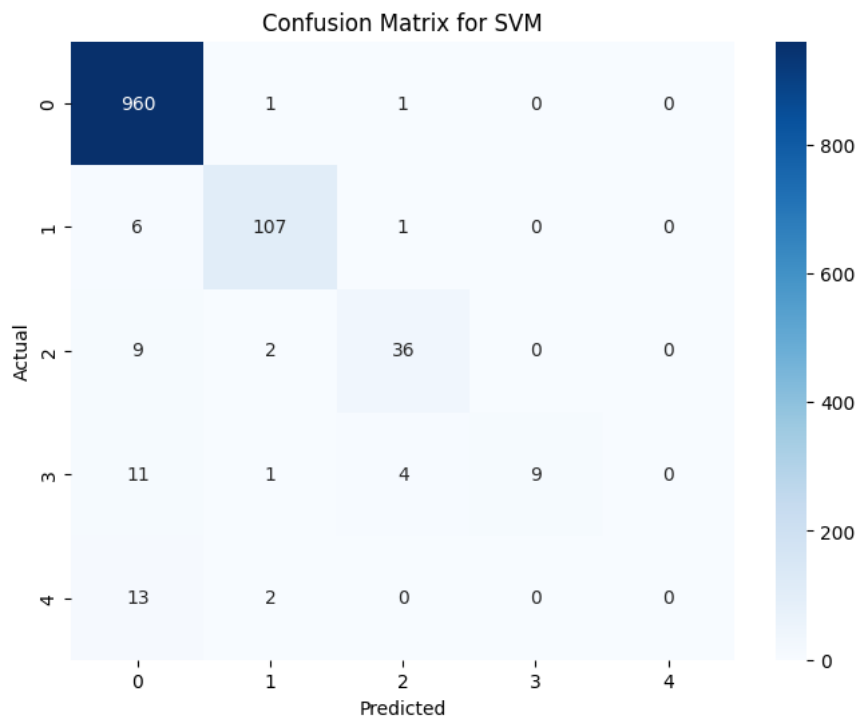


Figure 21: Confusion Matrix for SVM

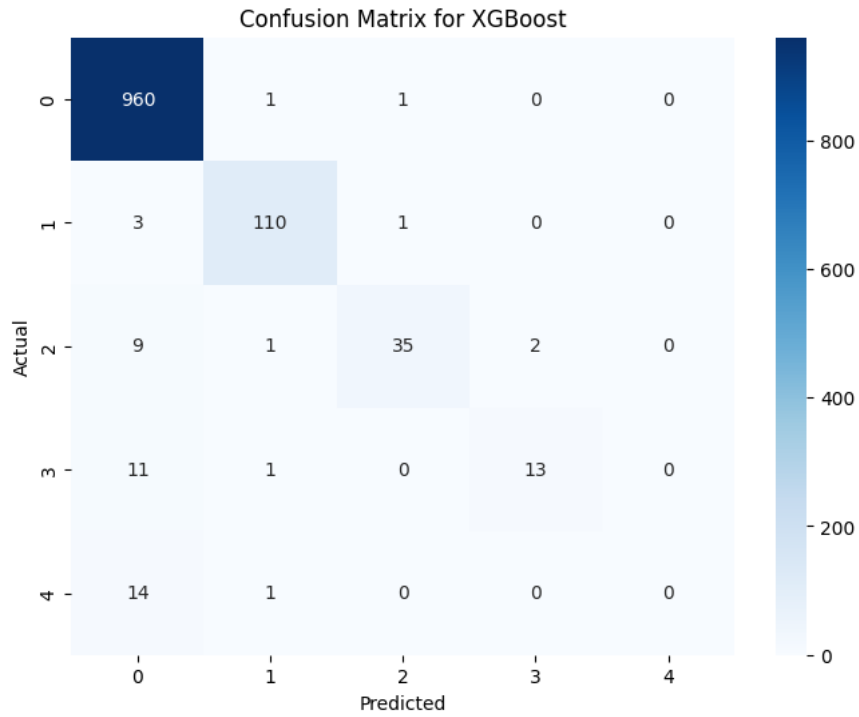


Figure 22: Confusion Matrix for XGBoost

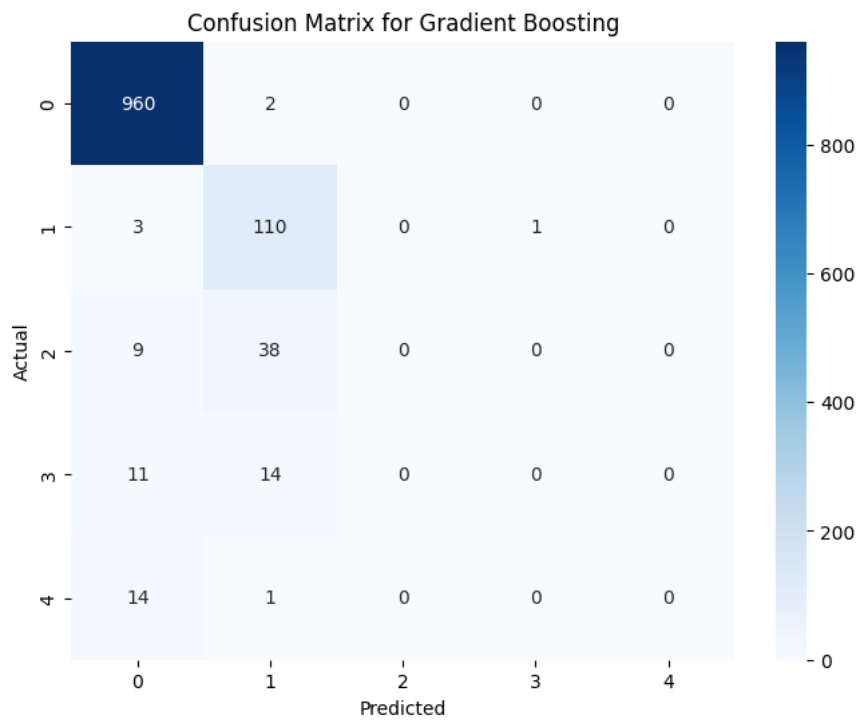


Figure 23: Confusion Matrix for Gradient Boosting

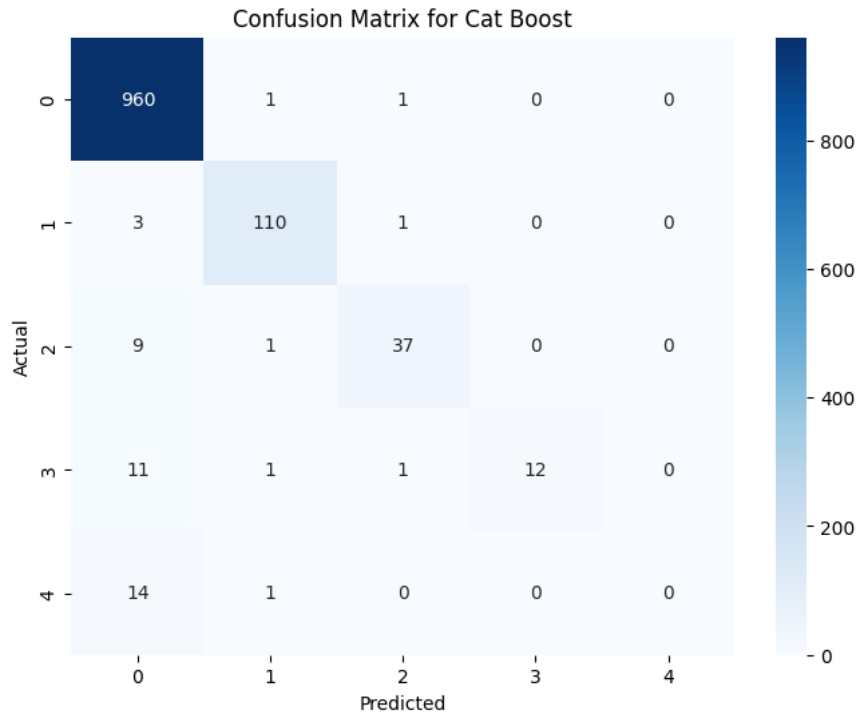


Figure 24: Confusion Matrix for Cat Boost

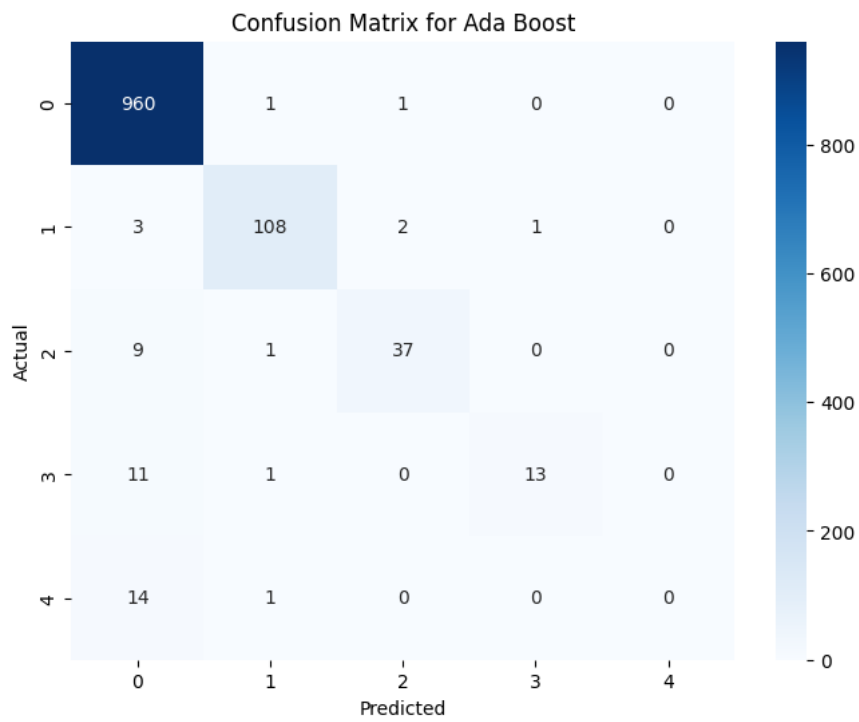


Figure 25: Confusion Matrix for Ada Boost

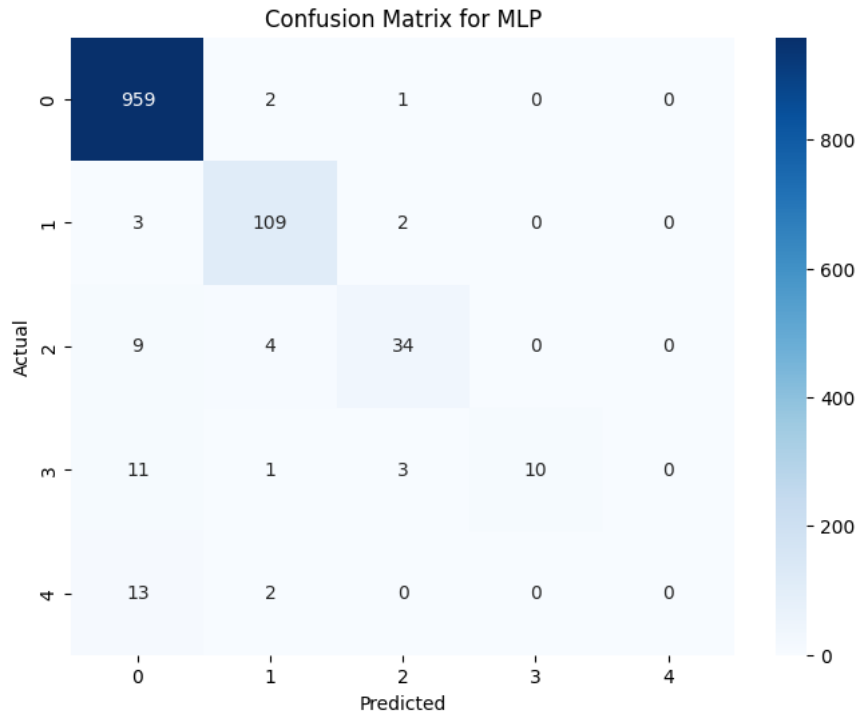


Figure 26: Confusion Matrix for MLP Classifier

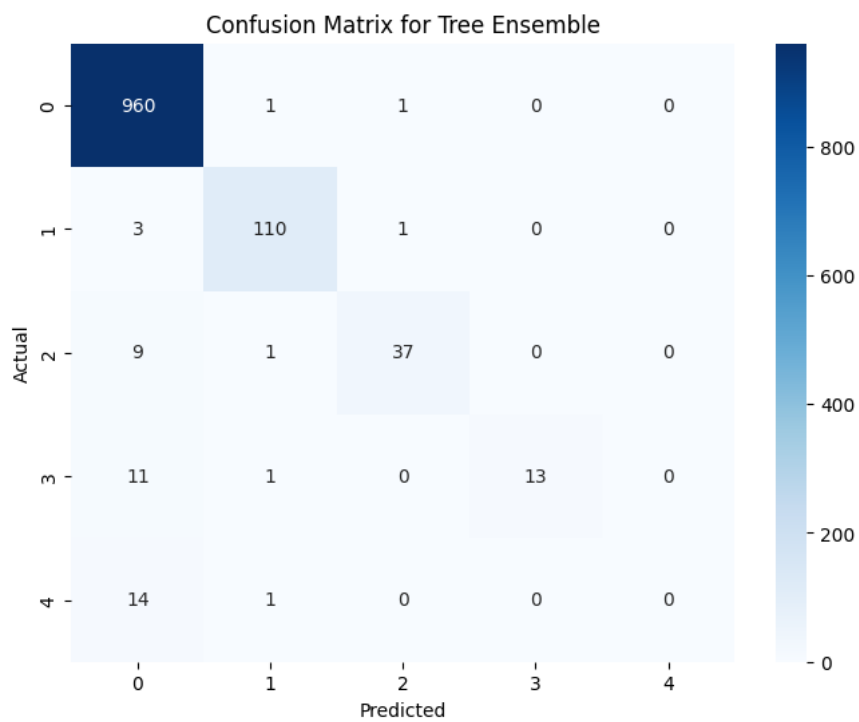


Figure 27: Confusion Matrix for Tree Ensemble

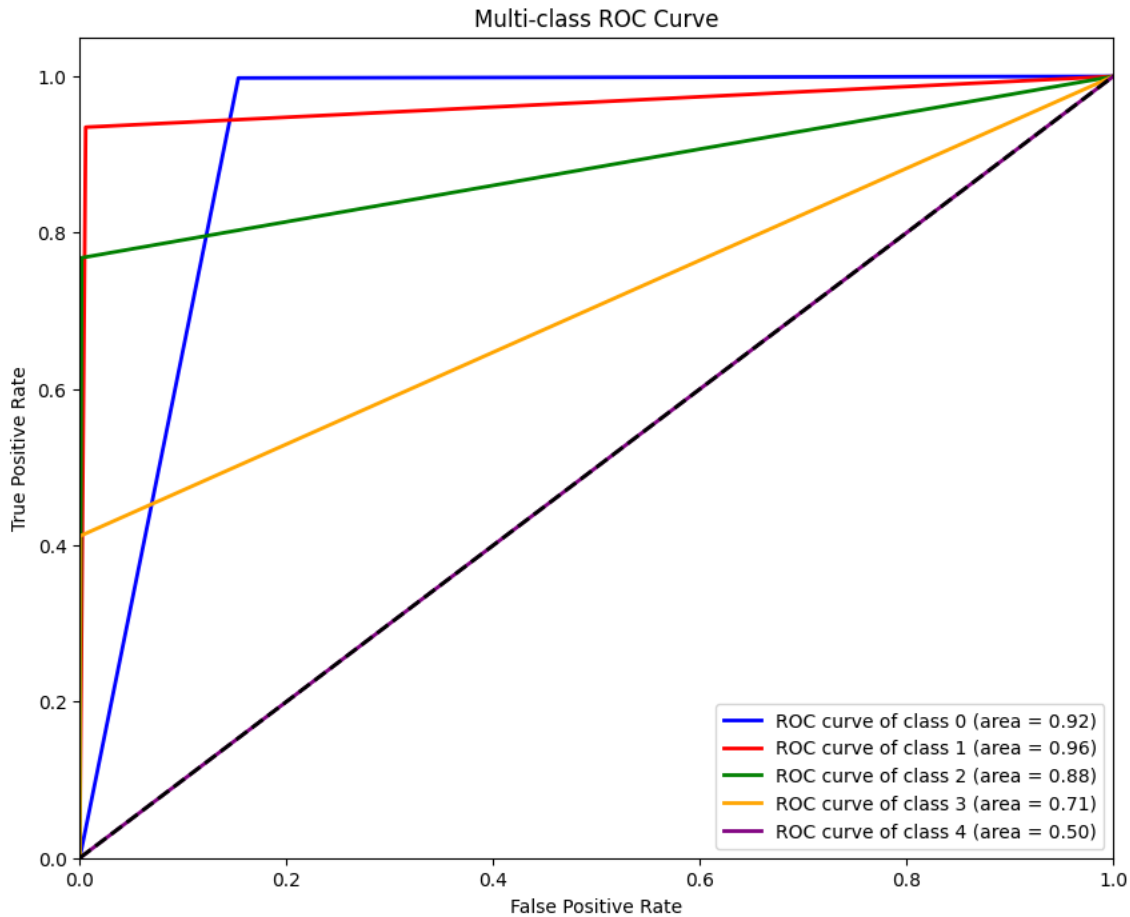


Figure 28: ROC Curve for Tree Ensemble

We can see the actual scenario of the result from the confusion matrices mention in Figure 16, 17,18,19,20,21,22,23,24,25,26,27 and 28.

4.3 Discussion

The comparison between the algorithms based on training and testing accuracies, as well as evaluation metrics like precision, recall, F1 score, and overall accuracy, reveals significant insights into their performance and suitability for the task.

Algorithms such as Decision Tree and Random Forest achieve perfect training accuracy (1.00) but show a slight drop in testing accuracy (0.93 and 0.95, respectively). This indicates potential overfitting, where these models perform exceptionally well on the training data but struggle to generalize to unseen data. On the other hand, ensemble methods such as Tree Ensemble, XGBoost, CatBoost, and AdaBoost maintain both high training (0.96) and testing accuracies (0.96–0.97), showcasing their ability to generalize effectively.

The metrics further highlight the strengths and weaknesses of each algorithm. Tree Ensemble and AdaBoost demonstrate the highest precision (0.77–0.78), recall (0.62), and F1

score (0.67), making them highly reliable for balancing true positive and false positive rates. Ensemble methods like Random Forest, XGBoost, and CatBoost also exhibit strong metrics, reflecting their robustness in handling complex patterns in the data. In contrast, simpler models like Naive Bayes and MLP Classifier exhibit lower precision (0.47), recall (0.44–0.47), and F1 scores (0.44–0.45), indicating weaker performance and reduced suitability for this dataset.

Tree Ensemble emerges as the best-performing model with the highest testing accuracy (0.97) and strong overall metrics, demonstrating excellent generalization and balanced predictions. Ensemble methods like XGBoost, CatBoost, and AdaBoost follow closely, offering competitive performance across all metrics.

The results emphasize the superiority of ensemble-based methods for achieving high accuracy and robust metric values, making them ideal for datasets with complex patterns. Simpler models like Naive Bayes struggle to match their performance, highlighting the importance of advanced techniques in modern machine learning tasks. The consistency of Tree Ensemble and other ensemble methods underscores their effectiveness in balancing precision, recall, and overall generalization.

The analysis highlights the strength of ensemble-based algorithms, particularly Tree Ensemble, XGBoost, and AdaBoost, in achieving high accuracy and balanced performance metrics. These algorithms demonstrate better generalization and robustness compared to simpler models like Naive Bayes and MLP Classifier, which show limitations in handling complex datasets. The results underscore the importance of advanced techniques in optimizing performance for real-world datasets.

CHAPTER 5

IMPACT ON SOCIETY, ENVIRONMENT, AND SUSTAINABILITY

5.1 Impact on Society

The application of machine learning approaches to assess the health impact of air pollution has profound implications for society, addressing critical environmental and public health challenges. These advancements offer transformative potential across various domains:

By providing accurate, real-time assessments of the health risks associated with air pollution, machine learning models enable early detection of high-risk scenarios. This facilitates timely interventions, reducing hospital admissions and mortality rates due to pollution-related diseases. Vulnerable populations, such as children, the elderly, and individuals with pre-existing conditions, particularly benefit from targeted strategies informed by these assessments.

Insights derived from machine learning models empower policymakers to design data-driven regulations and environmental standards aimed at reducing air pollution levels. Cities can leverage these findings to implement sustainable urban planning initiatives, such as optimizing traffic flow, enhancing green spaces, and regulating industrial emissions, fostering healthier living environments.

The reduction in pollution-related health issues decreases healthcare costs and minimizes productivity losses due to illness. Additionally, efficient resource allocation enabled by predictive models ensures that interventions, such as deploying air purifiers or relocating high-risk communities, are cost-effective and impactful.

The dissemination of localized pollution and health impact data raises public awareness, motivating individuals and communities to adopt eco-friendly practices. Increased societal understanding of pollution's consequences also drives demand for cleaner technologies and sustainable practices.

Machine learning approaches can address disparities in pollution monitoring and healthcare accessibility by offering scalable and cost-effective solutions for underserved regions. These tools help bridge gaps in knowledge and infrastructure, ensuring that low- and middle-income countries can effectively mitigate the health impacts of air pollution.

The integration of machine learning into environmental health research fosters interdisciplinary collaboration and accelerates innovation. These approaches pave the way for more

sophisticated models, datasets, and tools, contributing to the broader fields of artificial intelligence and public health.

By addressing the multifaceted challenges of air pollution and its health impacts, machine learning approaches hold the potential to significantly improve societal well-being, promote environmental sustainability, and create a healthier, more equitable future for all.

5.2 Impact on Environment

The application of machine learning approaches to assessing the health impact of air pollution has significant environmental implications, fostering more effective and sustainable solutions to address pollution challenges. Key impacts include:

Machine learning models enable the integration of data from diverse sources, such as satellite imagery, ground sensors, and meteorological data, providing real-time and accurate insights into pollution levels. This improves the ability to monitor and identify pollution hotspots, facilitating targeted mitigation strategies that reduce environmental degradation.

Insights derived from machine learning approaches empower policymakers to design informed and effective regulations to combat air pollution. By identifying critical pollution sources and trends, these models support the implementation of cleaner technologies, emission control measures, and sustainable practices that contribute to long-term environmental health.

Predictive analytics from machine learning can guide investments in renewable energy sources and eco-friendly technologies. For instance, optimizing energy consumption patterns and reducing industrial emissions are key outcomes that align with global sustainability goals, such as reducing greenhouse gas emissions.

By lowering pollutant levels, machine learning-informed interventions help protect ecosystems from harmful substances like sulfur dioxide (SO₂), nitrogen dioxide (NO₂), and particulate matter (PM_{2.5}). This aids in preserving biodiversity, maintaining soil fertility, and ensuring the health of aquatic ecosystems.

In urban areas, machine learning models can guide the development of green infrastructure, such as urban forests and green roofs, to combat pollution. These solutions not only improve air quality but also contribute to carbon sequestration, temperature regulation, and overall environmental resilience.

Reducing air pollution through data-driven policies and practices aligns with broader climate change mitigation strategies. Lower levels of pollutants like black carbon and methane, which contribute to global warming, directly impact the fight against climate change.

Machine learning approaches provide actionable insights that are accessible to both policymakers and the general public. By increasing awareness of pollution sources and their environmental impacts, these technologies encourage collective action and a shift toward environmentally responsible behaviors.

By fostering a deeper understanding of air pollution and its effects, machine learning not only mitigates immediate environmental challenges but also contributes to building a sustainable and healthier planet. These approaches ensure that interventions are both impactful and aligned with global environmental goals.

5.3 Ethical Aspects

The use of machine learning to assess the health impact of air pollution raises important ethical considerations that must be addressed to ensure responsible research and application. These aspects include:

Health data used in machine learning models often contains sensitive personal information. Ensuring that data is anonymized and securely stored is critical to protecting individual privacy and complying with regulations such as the General Data Protection Regulation (GDPR). Unauthorized access or misuse of data can lead to breaches of confidentiality and potential harm to individuals or communities. Machine learning models can inadvertently inherit biases from the datasets they are trained on, leading to unequal representation or inaccurate predictions for certain populations. Ensuring fairness in model development is essential to prevent the marginalization of vulnerable groups, such as low-income communities or regions with limited monitoring data. Complex machine learning models, especially deep learning approaches, often function as "black boxes," making it difficult to interpret their predictions. Transparent methodologies and clear explanations of model outputs are necessary to build trust among stakeholders, including policymakers, healthcare professionals, and the public.

The deployment of machine learning models should consider the accessibility of technology and resources in different regions, particularly in low- and middle-income countries. Failure to address resource disparities could exacerbate existing inequalities in environmental health management. Decisions informed by machine learning models, such as resource allocation or health interventions, must involve human oversight to ensure ethical judgment. Policymakers and practitioners must remain accountable for the outcomes of decisions based on model predictions. When using health data, individuals must provide informed consent, understanding how their data will be used, stored, and shared. Ensuring that

participants are aware of the purpose and scope of data collection fosters ethical data practices. Air pollution disproportionately affects marginalized communities. Ethical machine learning applications should aim to identify and address these disparities, supporting equitable interventions that prioritize the most affected populations. Ethical considerations must also address the potential long-term impacts of model deployment, such as unintended consequences on public policy or resource distribution. Continuous monitoring and adaptation of models are necessary to mitigate adverse outcomes. By addressing these ethical aspects, researchers and practitioners can ensure that machine learning approaches to assessing the health impact of air pollution are used responsibly and equitably, fostering trust and contributing to meaningful societal and environmental outcomes.

5.4 Sustainability Plan

To ensure the long-term viability and positive impact of using machine learning approaches to assess the health effects of air pollution, a comprehensive sustainability plan is essential. This plan addresses economic, social, environmental, and technological dimensions to create a robust framework for ongoing development and application.

1. Economic Sustainability:

- **Cost-Effective Data Collection:** Utilize scalable data sources such as satellite imagery, publicly available air quality sensors, and open health datasets to minimize operational costs.
- **Funding and Partnerships:** Secure long-term funding from governments, environmental organizations, and public health agencies while fostering partnerships with research institutions and private companies.
- **Open-Source Models and Tools:** Develop and share open-source machine learning models to reduce barriers to adoption and encourage collaborative improvements.

2. Environmental Sustainability:

- **Promoting Pollution Reduction Policies:** Use model insights to advocate for cleaner energy, sustainable urban planning, and industrial emission reductions.
- **Monitoring Environmental Progress:** Continuously evaluate the effectiveness of interventions by integrating real-time data and updating models to track progress toward air quality improvement goals.

- **Energy-Efficient Computing:** Implement energy-efficient algorithms and computational techniques to reduce the carbon footprint of model training and deployment.

3. Social Sustainability:

- **Equity and Accessibility:** Ensure that tools and insights are accessible to all regions, particularly underserved communities that face the greatest health risks from air pollution.
- **Capacity Building:** Train local governments, healthcare professionals, and environmental organizations to use and interpret machine learning outputs effectively.
- **Public Awareness and Engagement:** Disseminate findings in an understandable format to raise public awareness of air pollution's health impacts and encourage community action.

4. Technological Sustainability:

- **Scalable and Adaptable Models:** Develop models that can adapt to new data sources, emerging pollutants, and changing environmental conditions.
- **Cloud-Based Infrastructure:** Leverage cloud services to ensure scalable storage, computational power, and data sharing.
- **Interoperability:** Ensure that tools and models can integrate seamlessly with existing public health and environmental monitoring systems.

5. Governance and Policy Integration:

- **Stakeholder Collaboration:** Establish a multi-stakeholder governance model involving policymakers, researchers, and community leaders to ensure accountability and inclusivity.
- **Policy Support:** Provide actionable insights for evidence-based policymaking to promote sustainable urban development, energy transitions, and public health interventions.

6. Continuous Improvement:

- **Feedback Mechanisms:** Implement systems for continuous feedback from stakeholders to refine and improve the models.
- **Regular Updates:** Update machine learning models with new data, technologies, and methodologies to maintain accuracy and relevance.
- **Ethical Oversight:** Establish an ethical review board to monitor the fairness, transparency, and societal impact of the models.

By aligning economic, social, environmental, and technological efforts, this sustainability plan ensures that machine learning approaches to assessing the health impact of air pollution remain effective, equitable, and adaptable, fostering long-term societal and environmental benefits.

CHAPTER 6

CONCLUSION AND FUTURE WORK

6.1 Summary of the study

The study investigates the potential of machine learning techniques to evaluate the health impact of air pollution by leveraging diverse datasets and advanced algorithms. It identifies air pollution as a significant global health concern linked to various respiratory, cardiovascular, and chronic diseases, necessitating innovative approaches for real-time and localized assessments.

The research utilizes datasets containing air quality indices, pollutant measurements, meteorological factors, and health outcomes to develop predictive machine learning models. These models include ensemble-based methods (e.g., Random Forest, Gradient Boosting, XGBoost, and AdaBoost) and classical algorithms (e.g., Decision Tree, KNN, and SVM). The models are assessed using metrics such as precision, recall, F1 score, and accuracy.

The results highlight the effectiveness of ensemble methods, particularly Tree Ensemble, which achieves the highest accuracy and generalization. The study also underscores challenges like data heterogeneity, ethical considerations, and model interpretability.

The findings emphasize the transformative role of machine learning in monitoring pollution, predicting health impacts, and guiding policies. By integrating advanced analytics with environmental and health data, the study offers actionable insights for mitigating air pollution's adverse effects and promoting societal and environmental well-being.

6.2 Conclusions

This study demonstrates the transformative potential of machine learning in assessing the health impacts of air pollution, a pressing global concern linked to numerous chronic and acute health conditions. By integrating diverse datasets that encompass pollutant levels, meteorological factors, and health outcomes, advanced machine learning models were developed to predict and evaluate the health risks associated with air quality.

Among the algorithms evaluated, ensemble-based methods such as Tree Ensemble, XGBoost, and AdaBoost emerged as the most effective, achieving high accuracy and robust performance across multiple metrics. These models not only excelled in predictive capabilities but also proved their utility in capturing complex, non-linear relationships within the data.

The research also highlights key challenges, such as handling data heterogeneity, ensuring model interpretability, and addressing ethical concerns like fairness and privacy. Despite these challenges, the study underscores the potential of machine learning in providing actionable insights that can inform policy decisions, enable real-time monitoring, and ultimately improve public health outcomes.

In conclusion, the application of machine learning in this domain offers a promising pathway to mitigate the adverse health effects of air pollution. Future work should focus on refining these models by incorporating larger, more diverse datasets and enhancing model explainability to better support decision-makers and stakeholders in creating sustainable environmental and health policies.

6.3 Implication for further study

The findings of this study highlight several opportunities for advancing the application of machine learning in assessing the health impacts of air pollution. Integrating real-time data from IoT sensors and satellite imagery offers a pathway to develop dynamic models that provide predictions at finer temporal and spatial resolutions. This approach can enhance the applicability of these models in monitoring pollution hotspots and enabling timely interventions. Moreover, exploring advanced techniques like deep learning architectures, including RNNs and Transformer models, can help capture temporal dependencies and intricate relationships within the data, offering more robust and accurate predictions.

Future research can also focus on addressing data heterogeneity and enhancing model interpretability. Incorporating multi-domain datasets that include socio-economic factors, urban infrastructure data, and individual health profiles can provide a holistic understanding of pollution's impact on different populations. Additionally, the use of explainable AI (XAI) techniques, such as SHAP and LIME, can ensure model transparency and foster stakeholder trust, making machine learning outputs more accessible for decision-making. Ethical considerations, such as fairness and privacy, should also be prioritized to ensure equitable and responsible deployment of these technologies.

Lastly, further studies can expand on the integration of machine learning outputs with policy-making frameworks. By developing decision-support systems that combine predictive

models with health and economic impact assessments, researchers can guide policymakers in formulating effective interventions. Addressing the effects of emerging pollutants like microplastics and nanoparticles, as well as considering the influence of climate change on pollution dispersion, can further enhance the relevance of future studies. Promoting global collaboration and standardization of datasets and methodologies will ensure reproducibility and scalability, enabling more comprehensive and impactful research in this domain.

References

- [1] X. Yi, J. Zhang, Z. Wang, T. Li, and Y. Zheng, “Deep distributed fusion network for air quality prediction,” in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018.
- [2] Arvind Dabass *et al.*, “Association of exposure to particulate matter (PM_{2.5}) air pollution and biomarkers of cardiovascular disease risk in adult NHANES participants (2001–2008),” *International Journal of Hygiene and Environmental Health*, vol. 219, no. 3, pp. 301–310, May 2016, doi: <https://doi.org/10.1016/j.ijheh.2015.12.002>.
- [3] B. Zhai and J.-G. Chen, “Development of a stacked ensemble model for forecasting and analyzing daily average PM_{2.5} concentrations in Beijing, China,” *Science of The Total Environment*, vol. 635, pp. 644–658, Sep. 2018, doi: <https://doi.org/10.1016/j.scitotenv.2018.04.040>.
- [4] F. Xiao, M. Yang, H. Fan, G. Fan, and M. A. A. Al-qaness, “An improved deep learning model for predicting daily PM_{2.5} concentration,” *Scientific Reports*, vol. 10, no. 1, Dec. 2020, doi: <https://doi.org/10.1038/s41598-020-77757-w>.
- [5] T. V. Vu *et al.*, “Assessing the impact of clean air action on air quality trends in Beijing using a machine learning technique,” *Atmospheric Chemistry and Physics*, vol. 19, no. 17, pp. 11303–11314, Sep. 2019, doi: <https://doi.org/10.5194/acp-19-11303-2019>.
- [6] M. Castelli, F. M. Clemente, A. Popovič, S. Silva, and L. Vanneschi, “A Machine Learning Approach to Predict Air Quality in California,” *Complexity*, vol. 2020, pp. 1–23, Aug. 2020, doi: <https://doi.org/10.1155/2020/8049504>.
- [7] S. Ameer *et al.*, “Comparative Analysis of Machine Learning Techniques for Predicting Air Quality in Smart Cities,” *IEEE Access*, vol. 7, pp. 128325–128338, 2019, doi: <https://doi.org/10.1109/ACCESS.2019.2925082>.
- [8] Y.-C. Wu *et al.*, “Air quality monitoring using mobile microscopy and machine learning,” *Light: Science & Applications*, vol. 6, no. 9, pp. e17046–e17046, Sep. 2017, doi: <https://doi.org/10.1038/lsa.2017.46>.
- [9] W. Wei, O. Ramalho, L. Malingre, S. Sivanantham, J. C. Little, and C. Mandin, “Machine learning and statistical models for predicting indoor air quality,” *Indoor Air*, vol. 29, no. 5, pp. 704–726, Jul. 2019, doi: <https://doi.org/10.1111/ina.12580>.

- [10] D. Zhu, C. Cai, T. Yang, and X. Zhou, “A Machine Learning Approach for Air Quality Prediction: Model Regularization and Optimization,” *Big Data and Cognitive Computing*, vol. 2, no. 1, p. 5, Feb. 2018, doi: <https://doi.org/10.3390/bdcc2010005>.
- [11] Y.-C. Liang, Y. Maimury, A. H.-L. Chen, and J. R. C. Juarez, “Machine Learning-Based Prediction of Air Quality,” *Applied Sciences*, vol. 10, no. 24, p. 9151, Dec. 2020, doi: <https://doi.org/10.3390/app10249151>.
- [12] K. Meacham-Hensold *et al.*, “High-throughput field phenotyping using hyperspectral reflectance and partial least squares regression (PLSR) reveals genetic modifications to photosynthetic capacity,” *Remote Sens. Environ.*, vol. 231, no. 111176, p. 111176, 2019.
- [13] M. Jerrett *et al.*, “A review and evaluation of intraurban air pollution exposure models,” *J. Expo. Anal. Environ. Epidemiol.*, vol. 15, no. 2, pp. 185–204, 2005.
- [14] R. C. Nethery, F. Mealli, J. D. Sacks, and F. Dominici, “Causal inference and machine learning approaches for evaluation of the health impacts of large-scale air quality regulations,” *arXiv [stat.AP]*, 2019.
- [15] A. Mahajan, S. Mate, C. Kulkarni, and S. Sawant, “Predicting lung disease severity via image-based AQI analysis using deep learning techniques,” *arXiv [cs.CV]*, 2024.
- [16] X. Hu *et al.*, “Estimating PM_{2.5} concentrations in the conterminous United States using the random forest approach,” *Environ. Sci. Technol.*, vol. 51, no. 12, pp. 6936–6944, 2017.
- [17] *A Hybrid Model for PM_{2.5} Concentration Forecasting Based on Ensemble Empirical Mode Decomposition and Deep Learning*. .
- [18] Y. Zheng *et al.*, “Forecasting fine-grained air quality based on big data,” in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2015.
- [19] *Statistical Approaches to Address Multi-Pollutant Mixtures and Multiple Exposures: the ESCAPE Project*. .
- [20] *A Hybrid Forecasting Model for PM_{2.5} Mass Concentrations Based on Ensemble Empirical Mode Decomposition and a General Regression Neural Network*.
- [21] Rabie El Kharoua, “🌐 Air Quality and Health Impact Dataset🌐,” *Kaggle.com*, 2024. <https://www.kaggle.com/datasets/rabieelkharoua/air-quality-and-health-impact-dataset>.

Plagiarism Report:

Machine Learning Approaches to Assessing the Health Impact of Air Pollution

ORIGINALITY REPORT

12% SIMILARITY INDEX	13% INTERNET SOURCES	6% PUBLICATIONS	7% STUDENT PAPERS
--------------------------------	--------------------------------	---------------------------	-----------------------------

PRIMARY SOURCES

1	dspace.daffodilvarsity.edu.bd:8080 Internet Source	4%
2	Submitted to Daffodil International University Student Paper	2%
3	www.arxiv-vanity.com Internet Source	1%
4	www.catalyzex.com Internet Source	1%
5	tnsroindia.org.in Internet Source	1%
6	Submitted to University of Hong Kong Student Paper	1%
7	Submitted to CSU Northridge Student Paper	1%
8	Yi-Chen Wu, Ashutosh Shiledar, Yi-Cheng Li, Jeffrey Wong et al. "Air quality monitoring using mobile microscopy and machine learning", Light: Science & Applications, 2017 Publication	1%
9	acc-ern.tul.cz Internet Source	1%
10	Deepika Ghai, Kirti Rawal, Kanav Dhir, Suman Lata Tripathi. "Artificial Intelligence Techniques for Sustainable Development", CRC Press, 2024 Publication	1%