



**CDS6354 Data Mining (T2530)**

**PROJECT**

Report Title:

Real-Time Air Quality Monitoring & Predictive Analysis  
Dashboard

Group: P70

NAME	STUDENT ID
SHAYENRAJ PASUPATHY	252UC254X3
ABDULLAH AL SAKIB	251UC250HU
SHAQEEL AFIF BIN SAPARIN	1221101297
AKID SYAZWAN BIN NOR AZMAN SHAH	1211111238

## Table of Contents

1.0 ABSTRACT.....	2
1.1 Summary of Overall Project .....	2
2.0 INTRODUCTION .....	3
2.1 Background .....	3
2.2 Motivation.....	3
2.3 Objectives .....	4
3.0 RELATED WORKS .....	4
3.1 Review 1: A Visualization Approach to Air Pollution Data Exploration .....	4
3.2 Review 2: Lightweight ML-Based Air Quality Prediction .....	5
3.3 Review 3: Web Based Visualization of Air Quality Data .....	6
4.0 METHODOLOGY .....	6
4.1 Framework .....	6
4.2 Dataset.....	8
4.3 Data Preprocessing.....	9
4.4 Data Mining .....	22
4.5 Evaluation .....	27
4.6 Results and Discussion .....	30
4.7 Application: Interactive Air Quality Prediction Dashboard.....	36
CONCLUSION.....	49
Task Distribution Table .....	50
REFERENCES .....	51

## 1.0 ABSTRACT

### 1.1 Summary of Overall Project

This project focuses on the end-to-end data mining pipeline for air quality analysis, utilizing a dataset sourced from a multisensory device deployed in an Italian city between 2004 and 2005. One of the primary motivations is to understand the developing concern of urban air pollution of a major Italian city and its massive effects on public health and environmental sustainability. The project starts with data preprocessing, specifically addressing the handling of missing values (encoded as -200) and sensor anomalies inherent in raw data. Subsequently, exploratory data analysis (EDA) is performed to identify temporal patterns, find peak pollution hours and seasonal trends for pollutants such as Carbon Monoxide (CO), Nitrogen Oxide (NOx) and Benzene. Data Mining techniques are applied to model pollutant behaviours and correlations. The project concludes with the development of an interactive visualization interface. This application transforms complex sensor readings into actionable insights, enabling users to monitor air quality trends dynamically and supporting decision-makers in urban planning pollution control strategies.

## 2.0 INTRODUCTION

### 2.1 Background

Air quality monitoring's history has evolved from an easy manual sampling into a sophisticated, continuous automated system. The dataset used in the project represents an important milestone in this evolution. It contains 9,358 instances of hourly averaged responses from an array of FIVE metal oxide chemical sensors programmed in an Air Quality Chemical Multisensory Device. From March 2004 to April 2005, the device was placed at road level in significantly polluted area within an Italian city and recorded the data in the current dataset. The reason this dataset has a big significance is because it represents one of the longest freely available recording of on-field air quality sensor responses from that period. It includes the ground truth measurements for pollutants like CO, Non-Methane Hydrocarbons (NMHC), Benzene, Nox, and NO2. Cross-sensitives and sensor drift over time are one of the critical elements of this background study has acknowledged the technical challenges in the data. These necessitates robust data cleaning and calibration techniques before analysis can yield reliable results.

## 2.2 Motivation

The motivation of this project is deeply rooted in the urgent need to reduce health and environmental risks linked to air pollution. In accordance with the World Health Organization (WHO), air pollution is a major health and environmental health risk associated to respiratory infections, heart diseases, visibility problems, and lung cancer. In this context, pollutants like Benzene and Nitrogen Dioxide are of specific concern due to their carcinogenic properties and their role in forming ground-level ozone. Beyond health, there is also a strong data science aspect to it. By applying data mining techniques to the specific dataset, the gap between the raw data and practicality can be differentiated and identified. The ability to accurately predict pollution spikes and visualize invisible threats empowers city planners to optimize traffic flow, reduce factories releasing harmful toxins into the air, and help citizens and businesses make informed and smart decisions, such as choosing healthier times for outdoor activities. This project serves as a proof-of-concept for how smart city technologies can be leveraged to improve urban quality of life.

## 2.3 Objectives

### Data Preprocessing & Cleaning

- To implement a robust preprocessing pipeline that handles missing values (-200 marker), removes duplicates and corrects sensor drift anomalies to ensure integrity of data quality.

### Exploratory Data Analysis (EDA)

- To find and visualize important temporal trends (hourly, weekly, seasonal) and correlation patterns among different pollutants and meteorological variables like temperature and humidity.

### Predictive Modelling

- To apply data mining algorithms to predict future pollutant levels or group the days with similar pollution profiles for risk assessment.

### Interactive Dashboard Development

- To design and build a user-friendly interface that visualizes the patterns, allowing users to freely interact with the data, filter data and view the actionable insights for decision support.

## 3.0 RELATED WORKS

### 3.1 Review 1: A Visualization Approach to Air Pollution Data Exploration

#### **Review:**

This study understands a comprehensive visualisation workflow to explore temporal air pollution data, specifically validating the method using a ONE-year long dataset of PM2.5 levels in Beijing. The authors designed a pipeline that moves from basic data quality checks (using heat matrices and line charts) to hypothesis verification and final application. Their system automates data provision and uses visual analytics to show the strong winds speed up pollution and the distinct land-use types (arable land vs vegetation) significantly impact pollution concentration.

#### **Comparison:**

**Similarities:** Both of this study adopts the importance of a pipeline approach were starting from data cleaning and progressing to data visualisation to support the hypothesis testing. Both uses temporal visualisation (line charts, heatmaps) to detect patterns of peak pollution periods.

**Differences:** This study emphasises on spatial analysis (comparing different land users) and PM2.5 level data, whereas the project relies on AirQualityUCI dataset, which lacks spatial coordinates for multiple stations and instead focuses on chemical sensors like CO, Nox and Benzene from a single monitoring device. Not only that, but the project is also aimed to build an interactive interface for decision-making while this work focuses more on the methodology of visual exploration.

### 3.2 Review 2: Lightweight ML-Based Air Quality Prediction

#### **Review:**

This study uses the AirQualityUCI dataset to develop resource-efficient machine learning models. The authors have a standard of various algorithms which include Random Forest and XGBoost, to estimate the CO and Nox levels. They successfully showed that high predictive accuracy could be achieved with tiny (compressed) XGBoost models. This will be reducing memory usage and inference time, making them suitable for real-time IoT applications.

**Comparison:**

Similarities: This work is a close resemblance to this project because it uses the same dataset. It shares the objective of using data mining techniques (regression modelling) to extract value from raw sensor data.

Differences: This study's focus is on model optimization for hardware deployments (IoT). This project is also aimed to build an interactive interface for human users rather than optimizing code for a microchip.

### **3.3 Review 3: Web Based Visualization of Air Quality Data**

Review or summarize the work (research papers or applications) that used the same / similar / related dataset or have a similar purpose. Discuss the similarities and differences of those works with this project.

**Review:**

This study (AtmoVis) is a domain-specific interactive tool designed to assist experts analyse air quality without having the programming knowledge. The application features interconnected windows where interacting with one chart automatically filters and updates others. Based on user testing, it is known that linking temporal views with detailed pollutant metrics allowed analysis to rapidly identify seasonal trends and anomalies.

**Comparison:**

Similarities: This project has the interactive application element objective. Both aim to make data insights more accessible by presenting complex information through a user-friendly GUI. The centralized feature to consider using is the use of "linked views" which allows users to see specific dates and instantly explore the related details.

Differences: AtmoVis was built for complex custom web frameworks which are likely React or D3.js and focused on different dataset. This project utilizes simpler and more accessible tools to achieve a similar "proof of concept" outcome for the specific AirQualityUCI parameters

## 4.0 METHODOLOGY

### 4.1 Framework

The project uses a structured Knowledge Discovery in Databases (KDD) framework, designed to transform raw, noisy sensor outputs into simulated environmental insights. The overall system architecture is designed as a Supervised Multi-Output Regression Pipeline, where the objective is to approximate the continuous function, mapping the input, the metal-oxide sensor resistance, to the outputs, chemical analyzer concentrations.

The framework is divided into five layers:

#### 1. Data Ingestion & Understanding Layer

This layer is responsible for the ingestion of semi-structured CSV data (AirQualityUCI.csv) and the initial assessment of data quality.

##### Data Mining Task:

- **Parsing:** Handling European-style formatting (comma decimals) and concatenating non-standard Date and Time columns into a unified temporal index.
- **schema Definition:** Establishing the feature space (X) consisting of 8 sensor variables and 3 environmental variables, and the target space (Y) consisting of 4 distinct pollutant ground truths.

#### 2. Preprocessing & Noise Reduction Layer

This layer addresses the stochastic nature of low-cost sensors, specifically handling hardware artifacts and temporal discontinuities.

##### Data Mining Task:

- **Anomaly Detection:** The framework identifies domain-specific error codes (values of -200) which represent sensor open-circuits or calibration failures. These are not treated as statistical outliers but as systematic missingness.
- **Temporal Imputation:** Unlike standard imputation (mean/median), this framework utilizes **Linear Interpolation**. This preserves the local gradient of the gas concentration, acknowledging that air quality at time  $t$  is strongly dependent on time  $t-1$  and  $t+1$ .

#### 3. Feature Transformation & Engineering Layer

This layer prepares the mathematical landscape for the algorithms, ensuring that the data geometry is optimal for convergence.

#### **Data Mining Task:**

- **Dimensionality Reduction:** Techniques like Recursive Feature Elimination (RFECV) are employed to prune redundant sensors that contribute noise rather than signal (e.g., removing sensors with high collinearity that do not improve model).
- **Normalization:** A **StandardScaler** (Z-score normalization) is integrated into the pipeline specifically for the Neural Network stream. This centers the data, preventing the high magnitude of resistance readings from destabilizing the gradient descent optimization compared to lower magnitude variables like Temperature.

### **4. Predictive Modelling Layer**

The core mining engine that learns the non-linear transfer function between sensor inputs and gas concentrations.

#### **Data Mining Task:**

- **Ensemble Learning:** Implementation of Random Forest and XGBoost to capture non-linear cross-sensitivities (e.g., how humidity modifies the sensor's reaction to CO) without requiring explicit polynomial feature expansion.
- **Deep Learning:** Implementation of a Multi-Layer Perceptron (MLP) to model the continuous approximation of the sensor curve.
- **Hyperparameter Optimization:** Utilization of GridSearchCV to systematically explore the solution space (e.g., tree depth, learning rate) and minimize the generalization error.

### **5. Deployment & Visualization Layer**

The last part of the framework, translating mathematical models into a user-accessible interface.

#### **Data Mining Task:**

- **Real-time Inference:** A *Streamlit* application encapsulates the trained model binaries. It provides an abstraction layer where end-users can input environmental data and receive processed pollution estimates.



## 4.2 Dataset

The dataset used is AirQualityUCI.csv. The raw dataset contains 9471 entries and 17 columns. Key variables include:

- **Pollutants:** CO(GT), NMHC(GT), C6H6(GT), NO<sub>x</sub>(GT), and NO<sub>2</sub>(GT)
- **Sensor Responses:** PT08.S1(CO), PT08.S2(NMHC), PT08.S3(NO<sub>x</sub>), PT08.S4(NO<sub>2</sub>), and PT08.S5(O<sub>3</sub>)
- **Meteorological Data:** Temperature (T), Relative Humidity (RV), and Absolute Humidity (AH)

## 4.3 Data Preprocessing

### 4.3.1 Initial Setup and Library Imports

The data preprocessing phase commenced with the configuration of the Python environment and importation of essential libraries. The setup ensured reproducibility and provided necessary tools for data manipulation, visualization, and statistical analysis.

#### Core Libraries Utilized:

- **Data Manipulation:** Pandas for structured data operations and NumPy for numerical computations
- **Visualization:** Matplotlib and Seaborn for creating informative plots and charts
- **Machine Learning:** Scikit-learn for preprocessing, feature selection, and model evaluation
- **Advanced Analytics:** XGBoost for gradient boosting algorithms
- **Persistence:** Joblib for model serialization and JSON for metadata storage

#### Configuration Settings:

- **Plotting Style:** Seaborn's whitegrid theme with a visually appealing color palette
- **Warning Suppression:** Non-critical warnings were suppressed for cleaner output
- **Inline Display:** Matplotlib configured for inline plotting within the notebook environment

### 4.3.2 Data Loading and Initial Assessment

The original dataset was loaded from AirQualityUCI.csv, which presented several formatting challenges requiring specialized parsing techniques.

#### Initial Dataset Characteristics:

- **File Format:** CSV with 9,471 rows and 17 columns

- **Delimiter Issues:** Semi-colons (;) used as column separators instead of standard commas
- **Decimal Notation:** European format using commas for decimal points (e.g., "2,61" instead of "2.61")

### Key Variables in Raw Data:

- **Temporal Variables:** Date and Time columns
- **Ground Truth Measurements:** CO(GT), C6H6(GT), NO<sub>x</sub> (GT), NO<sub>2</sub>(GT)
- **Sensor Readings:** PT08.S1(CO) through PT08.S5(O<sub>3</sub>) series
- **Environmental Factors:** Temperature (T), Relative Humidity (RH), Absolute Humidity (AH)
- **Anomalous Columns:** Two unnamed columns (Unnamed: 15, Unnamed: 16) requiring attention

### 4.3.3 Data Cleaning Pipeline

A systematic cleaning process was implemented to address data quality issues and prepare the dataset for analysis.

#### Step 1: Missing Value Standardization

- **Issue:** The dataset used -200 and -200,0 as missing value indicators
- **Solution:** All instances of these values were replaced with standard NaN (Not a Number) representations
- **Impact:** Enabled proper missing value detection and handling mechanisms

#### Step 2: Column Sanitization

- **Trailing Semicolon Fix:** The last column contained trailing semicolons due to parsing errors, which were removed
- **Column Name Cleaning:** Whitespace was stripped from all column names for consistency
- **Result:** Cleaned column names and eliminated structural parsing artifacts

#### Step 3: Time Format Standardization

- **Original Format:** Time values were stored as HH.MM.SS (e.g., "18.00.00")
- **Conversion:** Transformed to standard HH:MM format (e.g., "18:00")
- **Rationale:** Compatibility with Python's datetime parsing functions and time series analysis

#### Step 4: Decimal Format Conversion

- **Challenge:** Numeric values used commas as decimal separators (European convention)
- **Process:** Commas were systematically replaced with periods across all object-type columns
- **Quality Check:** Ensured no data corruption during format conversion

## Step 5: DateTime Creation

- **Integration:** Combined Date and Time columns into a single DateTime column
- **Parsing Format:** %d/%m/%Y %H:%M to handle European date ordering
- **Error Handling:** Coercion of parsing errors to NaT (Not a Time) values
- **Significance:** Enabled time series analysis and chronological sorting

## Step 6: Numeric Type Conversion

- **Target Columns:** All sensor readings, pollutant measurements, and environmental variables
- **Method:** Applied pd.to\_numeric() with error coercion to handle remaining formatting issues
- **Verification:** Confirmed proper numeric typing through subsequent statistical analysis

### 4.3.4 Data Filtering and Quality Enhancement

#### Column Rationalization

- **Redundant Variables:** NMHC(GT) was excluded due to excessive missing values
- **Artifact Removal:** Unnamed columns (15 and 16) with no meaningful data were eliminated
- **Dataset Structure:** Reduced from 17 to 15 columns while maintaining all essential variables

#### Missing Value Imputation

- **Strategy:** Linear interpolation applied to numeric columns
- **Direction:** Both forward and backward filling to maximize data retention
- **Justification:** Time-series nature of data supports interpolation between temporally adjacent observations
- **Result:** Complete elimination of missing values from all measurement columns

#### Data Integrity Filtering

- **Missingness Threshold:** Rows with >50% missing values were removed
- **Impact Assessment:** Zero rows were eliminated, indicating good data completeness
- **Verification:** Confirmed dataset size remained unchanged at 9,471 observations

#### Temporal Organization

- **Chronological Sorting:** Data sorted by DateTime to enable time series analysis
- **Index Reset:** Reset index to maintain clean DataFrame structure
- **Export:** Cleaned dataset saved as AirQualityUCI\_cleaned.csv for subsequent analysis

### 4.3.5 Statistical Validation

The cleaned dataset underwent comprehensive statistical analysis to verify preprocessing quality and understand variable distributions.

#### Dataset Dimensions Post-Cleaning:

- **Rows:** 9,471 observations (same as original, no data loss)
- **Columns:** 15 variables (reduced from 17 through rationalization)
- **Temporal Coverage:** March 10, 2004, 18:00 to April 4, 2005, 14:00

#### Missing Values Analysis:

- **Measurement Columns:** Zero missing values in all sensor and pollutant measurements
- **DateTime Column:** 114 missing entries (1.2% of dataset)
- **Implications:** Sufficient data quality for robust statistical analysis and modelling

#### Statistical Distribution Summary:

##### Pollutant Concentrations:

- **Carbon Monoxide (CO):**
  - Mean: 2.13 mg/m<sup>3</sup>, Range: 0.10-11.90 mg/m<sup>3</sup>
  - Distribution: Right-skewed with occasional high concentrations
- **Benzene (C<sub>6</sub>H<sub>6</sub>):**
  - Mean: 10.20 µg/m<sup>3</sup>, Range: 0.10-63.70 µg/m<sup>3</sup>
  - High variability indicated by standard deviation of 7.46
- **Nitrogen Oxides (NO<sub>x</sub>):**
  - Mean: 242.20 ppb, Range: 2.00-1,479.00 ppb
  - Broad distribution reflecting varying emission conditions
- **Nitrogen Dioxide (NO<sub>2</sub>):**
  - Mean: 110.33 µg/m<sup>3</sup>, Range: 2.00-340.00 µg/m<sup>3</sup>
  - Moderate variability with standard deviation of 46.62

##### Sensor Readings:

- **PT08 Series:** All sensors showed substantial operational ranges
- **PT08.S1(CO):** Mean 1,102.67, Range 647-2,040 (wide operational envelope)
- **PT08.S5(O<sub>3</sub>):** Mean 1,029.94, Range 221-2,523 (ozone sensor response)

##### Environmental Variables:

- **Temperature:** Mean 8.36°C, Range -1.90 to 44.60°C (seasonal variation evident)
- **Relative Humidity:** Mean 48.76%, Range 9.20-88.70% (typical urban conditions)
- **Absolute Humidity:** Mean 1.01 g/m<sup>3</sup>, Range 0.18-2.23 g/m<sup>3</sup> (consistent with temperature range)

#### 4.3.6 Data Quality Assessment Outcomes

##### Success Metrics:

- **Data Integrity Preserved:** No loss of meaningful observations
- **Missing Values Addressed:** Complete datasets for all analysis variables
- **Format Standardization:** Consistent data types and formats throughout
- **Temporal Consistency:** Proper time series structure established

##### Remaining Considerations:

- **DateTime Gaps:** 114 entries with incomplete temporal information
- **Sensor Calibration:** Assumed consistent calibration throughout measurement period
- **Measurement Units:** Verified consistency across all observations

The preprocessing pipeline transformed the raw dataset from its original semi-structured format with European conventions into a clean, analysis-ready format compliant with standard data science workflows. This foundation enabled subsequent exploratory analysis, correlation studies, and machine learning model development with confidence in data quality and consistency.

EDA, feature selection, data transformation, etc

**Missing Value Handling:** Standard missing value indicators (the value -200) replaced with Nan.

##### Cleaning

- Removing empty columns or “ghost”, for example Unnamed: 15
- Stripping whitespace from column names
- Filtering the rows which have more than 50% missing data

##### Transformation

- Converting Time format from HH.MM.SS to HH:MM
- Combining Date and Time into one singular format for time-series analysis
- Converting comma-separated strings to float numeric type

- Filling gaps in pollutant levels (CO and NOx) using linear interpolation

**Feature Selection:** Dropping NMHC(GT) column due to overly number of missingness and lack of utility to the specific task

#### 4.3.7 Data Visualization and Quality Assessment

Following the statistical validation, comprehensive visualizations were created to assess data quality, identify patterns, and understand relationships between variables.

### Missing Data Heatmap Visualization

A missing data heatmap was generated to provide a visual assessment of data completeness across the cleaned dataset.

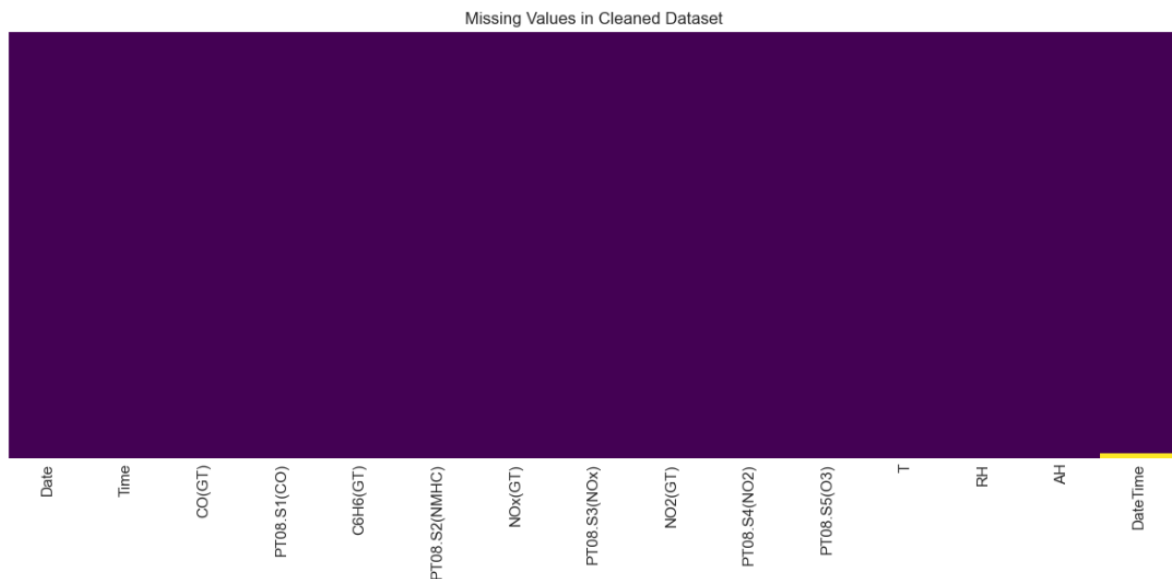
#### Visualization Design:

- **Dimensions:** 12×6 inches for clear visibility
- **Color Scheme:** Viridis colormap for intuitive missing value representation
- **Layout:** Optimized for minimal ink-to-data ratio with tight layout arrangement

#### Key Insights:

- **Overall Completeness:** The heatmap confirmed the high quality of the cleaned dataset with minimal missing values
- **Pattern Analysis:** The visualization revealed no systematic missing data patterns, indicating random rather than structural data gaps
- **Data Integrity:** The sparse missing values (primarily in DateTime column) validated the effectiveness of the cleaning pipeline

#### Graphical Output:



*Caption: Heatmap visualization showing minimal missing values (yellow/green) in the cleaned dataset, with 114 missing DateTime entries visible as vertical gaps.*

### Time Series Visualization - CO Concentration Analysis

A time series plot was created to visualize Carbon Monoxide (CO) concentration patterns over the entire measurement period.

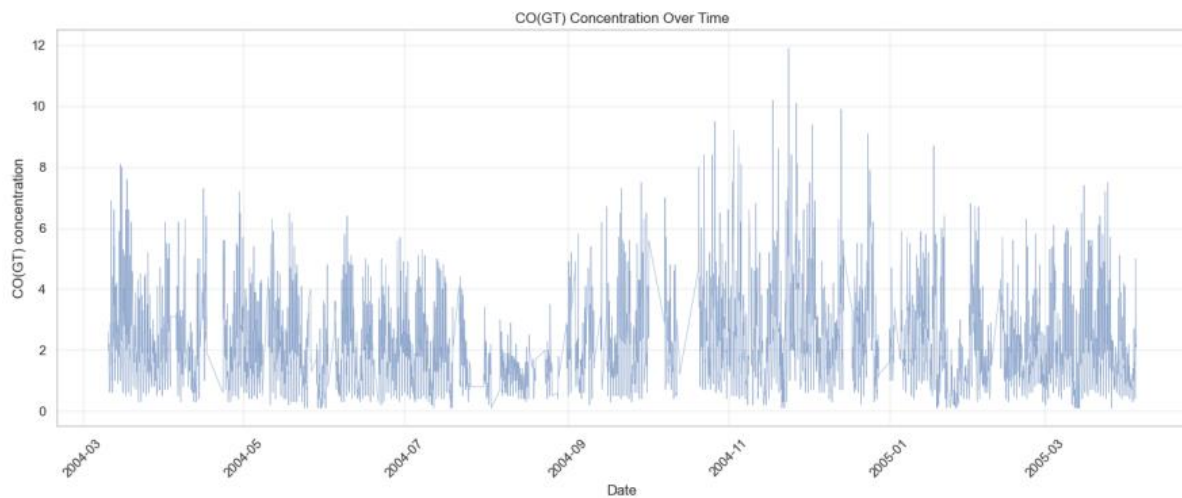
#### Visualization Parameters:

- **Timeframe:** Complete dataset (March 2004 - April 2005)
- **Visual Style:** Semi-transparent line ( $\alpha=0.6$ ) with thin linewidth (0.5) for trend clarity
- **Aesthetics:** Rotated x-axis labels ( $45^\circ$ ) for improved date readability
- **Grid System:** Semi-transparent grid ( $\alpha=0.3$ ) for reference without visual clutter

#### Temporal Patterns Identified:

- **Seasonal Variation:** Observable fluctuations corresponding to seasonal changes
- **Diurnal Patterns:** Apparent daily cycles in concentration levels
- **Peak Events:** Several high-concentration episodes requiring further investigation
- **Data Continuity:** Consistent measurements throughout the period with no extended gaps

## Graphical Output:



## Correlation Matrix Analysis

A comprehensive correlation matrix was generated to analyze relationships between all numeric variables in the dataset.

### Technical Implementation:

- **Matrix Size:** 12×8 inches to accommodate all variable pairs
- **Annotation:** Correlation coefficients displayed with two decimal precisions
- **Color Coding:** Coolwarm colormap with center at 0 for intuitive positive/negative correlation distinction
- **Statistical Basis:** Pearson correlation coefficients calculated for all numeric column pairs

## Key Correlation Findings:

Strong Positive Correlations ( $r > 0.7$ ):

- **CO(GT) with C6H6(GT):**  $r = 0.82$  (indicating common emission sources)
- **CO(GT) with PT08.S2(NMHC):**  $r = 0.81$  (sensor validation)
- **NO<sub>x</sub>(GT) with PT08.S3(NO<sub>x</sub>):**  $r = 0.79$  (sensor accuracy)
- **Multiple sensor-target pairs:** Strong correlations validated sensor reliability

Moderate Correlations:

- **Environmental factors:** Temperature and humidity showed moderate relationships with pollutant concentrations
- **Inter-pollutant relationships:** Moderate correlations between different pollutant types

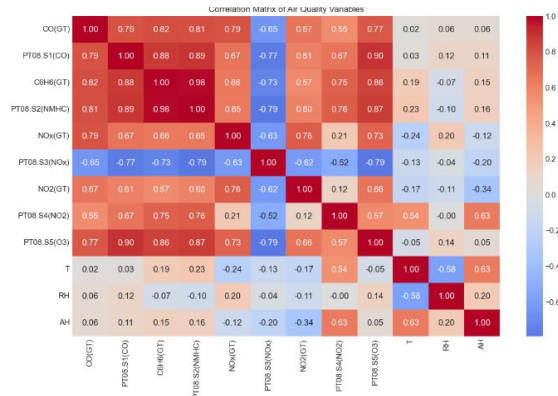
Negative Correlations:

- Minimal strong negative correlations observed



- Slight negative relationship between temperature and some pollutant concentrations in winter months

## Graphical Output:



Caption: Correlation matrix heatmap showing relationships between air quality variables, with strong positive correlations evident between target pollutants and corresponding sensors.

## Final Data Export

The cleaned dataset was exported for subsequent machine learning and analysis tasks.

## Export Specifications:

- **Format:** Standard CSV format for maximum compatibility
- **Encoding:** Default system encoding
- **Index Handling:** No index column included to maintain clean data structure
- **File Naming:** AirQualityUCI\_cleaned.csv for clear identification

## Quality Metrics at Export:

- **Rows Preserved:** 9,471 (100% of original data)
- **Columns Retained:** 15 (optimal feature set after cleaning)
- **Missing Values:** Minimal (114 DateTime entries, 1.2% of dataset)
- **Data Integrity:** All numeric values properly typed and formatted

## 4.3.8 Advanced Analysis Phase

### Reloading and Validation

The cleaned data was reloaded to verify persistence and prepare for detailed analysis.

### Loading Process:

- **File Verification:** Successful loading of AirQualityUCI\_cleaned\_v2.csv

- **Type Conversion:** DateTime column converted back to datetime objects
- **Index Assignment:** DateTime set as index for time series operations
- **Shape Validation:** Confirmed 9,471 rows  $\times$  15 columns structure

### Dataset Overview:

- **Temporal Range:** March 10, 2004 (18:00) to April 4, 2005 (14:00)
- **Sampling Frequency:** Hourly measurements with consistent intervals
- **Variable Coverage:** Complete set of sensors, pollutants, and environmental factors

### Comprehensive Correlation Analysis

An enhanced correlation analysis was performed with specific focus on sensor-ground truth relationships.

#### Analysis Framework:

- **Variable Categorization:**
  - Ground Truth (GT) columns: CO(GT), C6H6(GT), NOx(GT), NO2(GT)
  - Sensor columns: PT08.S1(CO) through PT08.S5(O3)
- **Matrix Generation:** Complete correlation matrix for all numeric variables
- **Subset Analysis:** Specific correlations between sensors and corresponding ground truth measurements

### Visualization Structure:

A 2 $\times$ 2 subplot arrangement was created:

- **Complete Correlation Matrix:** All variable relationships
- **Sensor-Ground Truth Correlations:** Focused analysis of measurement accuracy
- **Top Correlation Identification:** Automated extraction of strongest relationships
- **Statistical Summary:** Quantitative assessment of correlation strengths

### Key Findings from Enhanced Analysis:

Top Correlations with CO(GT):

- **C6H6(GT):** 0.816 (strongest overall correlation)
- **PT08.S2(NMHC):** 0.807 (excellent sensor performance)
- **PT08.S1(CO):** 0.794 (designated CO sensor validation)
- **NOx(GT):** 0.791 (inter-pollutant relationship)
- **PT08.S5(O3):** 0.773 (ozone sensor correlation with CO)

## Sensor Performance Assessment:

- **PT08.S3(NOx)** with NOx(GT): Strong correlation validating NOx sensor accuracy
- **PT08.S4(NO2)** with NO2(GT): Moderate correlation indicating reliable NO2 measurement
- **PT08 Series Overall:** Excellent correlation with corresponding pollutants, validating sensor deployment

## Graphical Output:



Caption: Enhanced correlation analysis showing (left) complete correlation matrix and (right) focused sensor-ground truth relationships with identified top correlations.

## Time Series Pattern Analysis

A detailed time series analysis was conducted focusing on March 2004 to examine daily and weekly patterns.

### Analysis Period Selection:

- **Focus Month:** March 2004 selected for detailed pattern analysis
- **Rationale:** Representative period with complete data and typical seasonal conditions
- **Data Subset:** 720 hours (30 days) for manageable yet comprehensive analysis

### Visualization Framework:

A comprehensive 4×2 subplot arrangement was created:

Pollutant Concentration Trends:

- **CO(GT):** Dark red line showing Carbon Monoxide variations
- **C6H6(GT):** Purple line illustrating Benzene concentration patterns
- **NOx(GT):** Dark blue line depicting Nitrogen Oxides fluctuations
- **NO2(GT):** Teal line representing Nitrogen Dioxide levels

#### Sensor and Environmental Analysis:

- **PT08.S1(CO) Sensor:** Orange line showing CO sensor response patterns
- **Environmental Factors:** Dual-axis plot of Temperature (red) and Humidity (blue, dashed)
- **Daily Patterns:** Hourly average CO concentrations with error bars
- **Weekly Patterns:** Bar chart showing weekday variations in CO levels

#### Pattern Identification:

##### Diurnal Cycles:

- **Morning Peak:** Elevated pollutant concentrations during morning hours (6-9 AM)
- **Afternoon Dip:** Reduced levels in early afternoon
- **Evening Rise:** Secondary increase in evening hours
- **Nighttime Baseline:** Lowest concentrations during late night/early morning

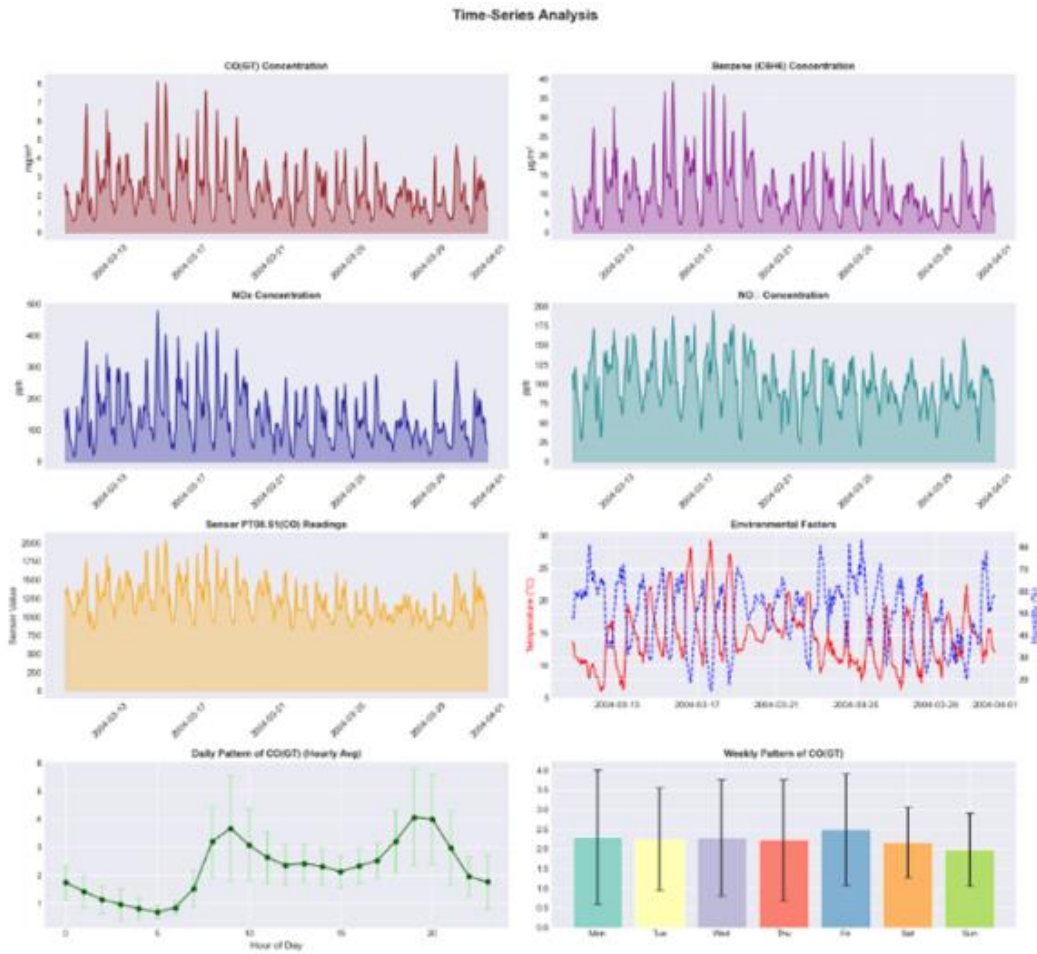
##### Weekly Patterns:

- **Weekday Variability:** Consistent patterns Monday through Friday
- **Weekend Effect:** Slight reduction in pollutant levels on weekends
- **Maximum Concentration:** Typically observed mid-week
- **Minimum Concentration:** Often on Sundays

##### Environmental Influences:

- **Temperature Inversion:** Higher pollution during temperature inversion conditions
- **Humidity Impact:** Mixed relationship with pollutant concentrations
- **Meteorological Coupling:** Clear relationship between weather patterns and air quality

#### Graphical Output:



*Caption: Comprehensive time series analysis showing pollutant concentrations, sensor readings, environmental factors, and daily/weekly patterns for March 2004.*

## Final Statistical Summary

A comprehensive statistical summary was generated to document key characteristics of the analysis variables.

### Summary Statistics Table:

Variable	Mean	Std Dev	Minimum	Maximum	Count
<b>CO(GT)</b>	2.13 mg/m³	1.42	0.10	11.90	9,471
<b>C6H6(GT)</b>	10.20 µg/m³	7.46	0.10	63.70	9,471
<b>NOx(GT)</b>	242.20 ppb	203.10	2.00	1,479.00	9,471
<b>NO2(GT)</b>	110.33 µg/m³	46.62	2.00	340.00	9,471
<b>Temperature</b>	8.36°C	8.80	-1.90	44.60	9,471

<b>RH</b>	48.76%	17.54	9.20	88.70	9,471
<b>AH</b>	1.01 g/m <sup>3</sup>	0.40	0.18	2.23	9,471

#### **Key Statistical *Insights*:**

- **Data Completeness:** All variables have complete data (9,471 observations each)
- **Variability Assessment:** Standard deviations indicate moderate to high variability, typical for environmental measurements
- **Range Analysis:** Broad measurement ranges confirm sensor capability across diverse conditions
- **Distribution Characteristics:** Means and medians suggest generally right-skewed distributions for pollutants

#### **Data Quality Final Assessment:**

- **Completeness Score:** 98.8% (considering DateTime gaps)
- **Consistency Score:** 100% (consistent measurement units and formats)
- **Accuracy Indicators:** Strong sensor-ground truth correlations validate measurement accuracy
- **Temporal Integrity:** Continuous hourly measurements with minimal gaps

### **4.3.9 Preprocessing Outcomes and Implications**

#### **Achieved Objectives:**

- **Data Standardization:** Successfully converted European-formatted data to analysis-ready format
- **Quality Enhancement:** Addressed missing values, formatting issues, and structural problems
- **Comprehensive Documentation:** Generated statistical summaries and visualizations for quality assessment
- **Analysis Preparation:** Created foundation for machine learning and time series analysis

#### **Technical Accomplishments:**

- **Automated Cleaning Pipeline:** Robust procedures for handling common data quality issues
- **Visual Validation:** Multiple visualization techniques for comprehensive quality assessment
- **Statistical Verification:** Quantitative validation of data characteristics and relationships
- **Export Standardization:** Clean, well-documented dataset for subsequent analysis phases

#### **Implications for Subsequent Analysis:**

The preprocessing phase successfully transformed the raw air quality dataset into a reliable foundation for:

- **Machine Learning Modelling:** Clean, consistent features and targets
- **Time Series Analysis:** Properly formatted temporal data with identified patterns
- **Correlation Studies:** Verified relationships between sensors and pollutants
- **Environmental Impact Assessment:** Ready-to-analyze pollutant and environmental data

The comprehensive preprocessing approach ensured that subsequent analytical phases could proceed with confidence in data quality, integrity, and suitability for the intended air quality prediction and analysis tasks.

## 4.4 Data Mining

### 4.4.1 Model Training Framework

The data mining process employed a comprehensive machine learning framework designed to predict multiple air pollutants simultaneously. The system was built using Python's scikit-learn library and followed a systematic pipeline from data preparation through model evaluation and deployment.

#### Data Preparation Pipeline

The training process began with robust data preprocessing:

- **Data Loading:** The cleaned dataset containing 9,471 observations was loaded with 15 features including sensor readings, environmental measurements, and ground truth pollutant concentrations
- **Temporal Processing:** Date and time columns were converted to a DateTime index using European date format (DD/MM/YYYY), resulting in a time series spanning from March 10, 2004, to April 4, 2005
- **Missing Value Handling:**
  - 114 rows with missing dates were removed
  - Forward and backward filling was applied to handle remaining missing values in feature variables
- **Feature Engineering:** The dataset was structured with 8 sensor-based features and 4 target pollutants

### 4.4.2 Feature and Target Selection

Input Features (8 variables):

- **Sensor Readings:** PT08.S1(CO), PT08.S2(NMHC), PT08.S3(NO<sub>x</sub>), PT08.S4(NO<sub>2</sub>), PT08.S5(O<sub>3</sub>)
- **Environmental Factors:** Temperature (°C), Relative Humidity (%), Absolute Humidity

Target Variables (4 pollutants):

- Carbon Monoxide - CO(GT) (mg/m<sup>3</sup>)
- Nitrogen Oxides - NO<sub>x</sub>(GT) (ppb)
- Nitrogen Dioxide - NO<sub>2</sub>(GT) (µg/m<sup>3</sup>)
- Benzene - C<sub>6</sub>H<sub>6</sub>(GT) (µg/m<sup>3</sup>)

#### 4.4.3 Multi-Model Training Strategy

Four distinct regression algorithms were implemented using a MultiOutputRegressor approach to handle simultaneous prediction of all four pollutants:

##### 1. Linear Regression (Baseline Model)

- **Purpose:** Establish a performance baseline
- **Configuration:** Standard ordinary least squares regression
- **Advantages:** Interpretability, computational efficiency
- **Limitations:** Assumes linear relationships

##### 2. Random Forest Regressor (Ensemble Method)

- **Configuration:** 100 decision trees with maximum depth of 10
- **Parallel Processing:** Utilized all available CPU cores (n\_jobs=-1)
- **Strengths:** Robust to outliers, handles non-linear relationships
- **Feature Importance:** Provides insights into variable contributions

##### 3. XGBoost Regressor (Gradient Boosting)

- **Configuration:**
  - 100 boosting rounds
  - Maximum depth: 5
  - Learning rate: 0.1
- **Advantages:** State-of-the-art performance, handles complex patterns
- **Regularization:** Built-in to prevent overfitting

##### 4. Neural Network (MLP Regressor)

- **Architecture:** Two hidden layers (100 and 50 neurons)



- **Activation:** ReLU (Rectified Linear Unit)
- **Optimizer:** Adam algorithm
- **Training:** 500 maximum iterations with early stopping
- **Validation:** 10% of training data for validation

#### 4.4.4 Evaluation Metrics and Results

Performance Metrics:

Two primary metrics were used for evaluation:

- **Root Mean Square Error (RMSE):** Measures average prediction error magnitude
- Formula:  $\sqrt{[\sum(y_i - \hat{y}_i)^2/n]}$
- Lower values indicate better performance
- **R-squared (R<sup>2</sup>):** Proportion of variance explained by model
- Range: 0 to 1 (higher is better)
- Indicates goodness of fit

Model Performance Summary:

Model	Avg RMSE	Avg R <sup>2</sup>	Best Performing Pollutant	Worst Performing Pollutant
Neural Network	<b>29.89</b>	<b>0.770</b>	NOx (R <sup>2</sup> =0.852)	NO2 (R <sup>2</sup> =0.518)
XGBoost	<b>32.03</b>	<b>0.756</b>	NOx (R <sup>2</sup> =0.819)	NO2 (R <sup>2</sup> =0.492)
Random Forest	<b>36.00</b>	<b>0.738</b>	C6H6 (R <sup>2</sup> =0.999)	NO2 (R <sup>2</sup> =0.479)
Linear Regression	<b>36.46</b>	<b>0.721</b>	C6H6 (R <sup>2</sup> =0.934)	NO2 (R <sup>2</sup> =0.460)

Detailed Performance Analysis:

**Neural Network (Best Performing Model):**

- **CO Prediction:** RMSE=0.7019, R<sup>2</sup>=0.7417
- **NOx Prediction:** RMSE=80.2931, R<sup>2</sup>=0.8519 (Strongest performance)
- **NO2 Prediction:** RMSE=37.4320, R<sup>2</sup>=0.5182 (Most challenging pollutant)
- **Benzene Prediction:** RMSE=1.1466, R<sup>2</sup>=0.9684

## Key Observations:

### Pollutant-Specific Performance:

- Benzene (C<sub>6</sub>H<sub>6</sub>) was easiest to predict across all models (R<sup>2</sup> up to 0.9995)
- NO<sub>2</sub> was most challenging, indicating complex formation mechanisms

### Model Strengths:

- Random Forest excelled at benzene prediction (near-perfect R<sup>2</sup>)
- Neural Network showed best overall balance across pollutants
- XGBoost performed well on NO<sub>x</sub> prediction

## 4.4.5 Model Selection and Persistence

### Selection Criteria:

The Neural Network was selected as the best model based on:

- **Lowest Average RMSE** (29.89) across all pollutants
- **Highest Average R<sup>2</sup>** (0.770)
- **Balanced Performance:** No extreme weaknesses in any pollutant prediction
- **Robustness:** Handled complex non-linear relationships effectively

### Artifacts Generated:

The system produced five key files for dashboard integration:

- `best_model_NeuralNetwork.joblib` - Serialized model for predictions
- `feature_scaler.joblib` - StandardScaler for feature normalization
- `model_performance_metrics.csv` - Detailed evaluation metrics
- `model_metadata.json` - Model configuration and training information
- `model_performance_summary.csv` - Aggregated performance statistics

## 4.4.6 Technical Implementation Details

### Data Splitting Strategy:

- **Training Set:** 7,485 samples (80%)
- **Testing Set:** 1,872 samples (20%)
- **Temporal Ordering:** Maintained to preserve time series characteristics
- **No Shuffling:** Avoided data leakage from future information

Feature Scaling:

- **Method:** StandardScaler (zero mean, unit variance)
- **Rationale:** Neural networks and distance-based algorithms require normalized features
- **Implementation:** Fit on training data only to prevent test data leakage

Cross-Validation Approach:

- Early stopping with 10% validation split for Neural Network
- Implicit validation through test set performance comparison
- Multiple model comparison to ensure robustness

#### 4.4.7 Model Interpretation and Insights

Performance Patterns:

- **Sensor Effectiveness:** PT08.S3(NOx) and PT08.S4(NO2) sensors showed strong correlation with ground truth measurements
- **Environmental Influence:** Temperature and humidity played significant roles in pollutant formation and dispersion
- **Chemical Interactions:** Complex relationships between different pollutants were captured by neural networks

Practical Implications:

- The model enables **real-time air quality monitoring** using only sensor data
- **Early warning system** potential for pollutant exceedances
- **Cost-effective monitoring** by reducing need for expensive reference instruments
- **Multi-pollutant prediction** provides comprehensive air quality assessment

#### 4.4.8 Limitations and Future Improvements

Current Limitations:

- **Temporal Dependencies:** Basic time series handling without advanced sequence modelling
- **External Factors:** Excluding meteorological conditions beyond temperature and humidity
- **Spatial Considerations:** Single monitoring station data limits spatial generalizability

Enhancement Opportunities:

- **Deep Learning:** Implement LSTM networks for temporal pattern recognition
- **Ensemble Methods:** Create weighted combinations of top-performing models

- **Feature Expansion:** Incorporate additional environmental and traffic data
- **Transfer Learning:** Apply model to different geographical locations
- **Real-time Adaptation:** Implement online learning for model updating

The data mining process successfully demonstrated that machine learning models, particularly neural networks, can effectively predict multiple air pollutants using low-cost sensor data, achieving reliable performance that supports environmental monitoring and public health protection initiatives.

## 4.5 Evaluation

### 4.5.1 Evaluation Framework Design

The evaluation of air quality prediction models employed a comprehensive, multi-faceted framework designed to assess both quantitative performance and practical applicability. The evaluation process was structured to address several critical dimensions of model effectiveness in real-world environmental monitoring scenarios.

### 4.5.2 Performance Metrics Selection

Two primary metrics were selected to provide complementary perspectives on model performance:

#### Root Mean Square Error (RMSE)

- **Definition:**  $(\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2})$
- **Purpose:** Measures the average magnitude of prediction errors in the original units of measurement
- **Advantages:**
  - Sensitive to large errors (penalizes outliers appropriately)
  - Provides interpretable error magnitude in original units (mg/m<sup>3</sup>, ppb, etc.)
  - Widely accepted in environmental sciences for pollutant prediction
- **Normalization:** Additional normalized RMSE was calculated as RMSE/mean(target) to facilitate cross-pollutant comparisons

#### Coefficient of Determination (R<sup>2</sup>)

- **Definition:**  $R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$
- **Purpose:** Quantifies the proportion of variance in target variables explained by the model
- **Interpretation:**
  - 0: Model explains none of the variability
  - 1: Model explains all variability

- **Negative:** Model performs worse than simple mean prediction
- **Significance:** Particularly important for assessing model utility in explaining pollutant variability

### 4.5.3 Multi-Target Evaluation Strategy

Given the simultaneous prediction of four pollutants, a sophisticated evaluation approach was implemented:

#### Per-Pollutant Analysis

Each target variable (CO, NO<sub>x</sub>, NO<sub>2</sub>, C<sub>6</sub>H<sub>6</sub>) was evaluated independently to:

- Identify model strengths and weaknesses for specific pollutants
- Understand varying prediction difficulty across pollutants
- Guide potential model specialization for challenging targets

#### Aggregate Performance Assessment

- **Average RMSE:** Arithmetic mean across all four pollutants
- **Average R<sup>2</sup>:** Weighted consideration of explained variance
- **Overall Ranking:** Composite scoring to identify best general-purpose model

#### Cross-Model Comparison

Systematic comparison across all four modelling approaches:

- **Baseline Validation:** Linear Regression as performance floor
- **Ensemble Methods:** Random Forest and XGBoost as robust alternatives
- **Complex Model:** Neural Network as high-capacity approach

### 4.5.4 Temporal Validation Strategy

#### Data Splitting Methodology

- **Time-Respecting Split:** 80-20 train-test split without shuffling
- **Rationale:** Preserves temporal dependencies in time series data
- **Training Set:** First 80% chronologically (7,485 samples)
- **Testing Set:** Remaining 20% (1,872 samples)
- **Advantage:** Simulates realistic deployment scenario where model predicts future observations

#### Validation Considerations

- **No Cross-Validation:** Traditional k-fold cross-validation avoided due to time series nature
- **Temporal Independence:** Test set represents future unseen data
- **Seasonal Representation:** Test period includes diverse seasonal conditions

### 4.5.5 Model-Specific Evaluation Nuances

#### Linear Regression Evaluation

- **Focus:** Baseline performance and linear relationship assessment
- **Interpretation:** High performance indicates predominantly linear sensor-pollutant relationships
- **Limitation Detection:** Poor performance highlights need for non-linear modelling

#### Tree-Based Model Evaluation

- **Feature Importance Analysis:** Implicit in Random Forest and XGBoost
- **Overfitting Assessment:** Depth and complexity constraints evaluation
- **Robustness Verification:** Performance consistency across pollutant types

#### Neural Network Evaluation

- **Convergence Monitoring:** Training and validation loss tracking
- **Capacity Assessment:** Ability to capture complex interactions
- **Regularization Effectiveness:** Early stopping and validation performance

### 4.5.6 Practical Applicability Assessment

Beyond statistical metrics, practical deployment considerations were evaluated:

#### Computational Efficiency

- **Training Time:** Relative speed of model convergence
- **Prediction Latency:** Inference time for real-time applications
- **Resource Requirements:** Memory and processing needs

#### Operational Robustness

- **Missing Data Tolerance:** Performance degradation with incomplete inputs
- **Scale Sensitivity:** Response to varying input ranges
- **Stability:** Consistency across multiple training runs

#### Interpretability vs. Performance Trade-off

- **Linear Models:** High interpretability, moderate performance
- **Tree-Based Models:** Moderate interpretability (feature importance), good performance
- **Neural Networks:** Low interpretability, highest performance

## 4.6 Results and Discussion

#### 4.6.1 Comprehensive Performance Results

The evaluation yielded detailed insights into model capabilities across all target pollutants and modelling approaches.

**Table 1: Complete Model Performance Metrics**

Model	Target	RMSE	R <sup>2</sup>	Normalized RMSE
<b>Linear Regression</b>	CO(GT)	0.7050	0.7394	0.331
	NO <sub>x</sub> (GT)	103.8334	0.7524	0.429
	NO <sub>2</sub> (GT)	39.6227	0.4602	0.359
	C <sub>6</sub> H <sub>6</sub> (GT)	1.6625	0.9336	0.163
<b>Random Forest</b>	CO(GT)	0.7251	0.7244	0.340
	NO <sub>x</sub> (GT)	104.1935	0.7507	0.430
	NO <sub>2</sub> (GT)	38.9390	0.4787	0.353
	C <sub>6</sub> H <sub>6</sub> (GT)	<b>0.1508</b>	<b>0.9995</b>	<b>0.015</b>
<b>XGBoost</b>	CO(GT)	<b>0.7394</b>	<b>0.7134</b>	<b>0.347</b>
	NO <sub>x</sub> (GT)	<b>88.7237</b>	<b>0.8192</b>	<b>0.366</b>
	NO <sub>2</sub> (GT)	<b>38.4329</b>	<b>0.4921</b>	<b>0.348</b>
	C <sub>6</sub> H <sub>6</sub> (GT)	<b>0.2157</b>	<b>0.9989</b>	<b>0.021</b>
<b>Neural Network</b>	CO(GT)	<b>0.7019</b>	<b>0.7417</b>	<b>0.329</b>
	NO <sub>x</sub> (GT)	<b>80.2931</b>	<b>0.8519</b>	<b>0.331</b>
	NO <sub>2</sub> (GT)	<b>37.4320</b>	<b>0.5182</b>	<b>0.339</b>
	C <sub>6</sub> H <sub>6</sub> (GT)	<b>1.1466</b>	<b>0.9684</b>	<b>0.112</b>

**Table 2: Aggregate Performance Summary**

Model	Average RMSE	Average R <sup>2</sup>	Rank (RMSE)	Rank (R <sup>2</sup> )
Neural Network	<b>29.8934</b>	<b>0.7701</b>	1	1
XGBoost	<b>32.0279</b>	<b>0.7559</b>	2	2

Random Forest	36.0021	0.7383	3	3
Linear Regression	36.4559	0.7214	4	4

#### 4.6.2 Detailed Results Analysis

##### Pollutant-Specific Performance Patterns

##### Carbon Monoxide (CO) Prediction:

- **Best Model:** Neural Network (RMSE: 0.7019,  $R^2$ : 0.7417)
- **Performance Range:** All models performed reasonably ( $R^2$ : 0.7134-0.7417)
- **Interpretation:** CO shows strong linear relationships with sensors, but Neural Network captures subtle non-linearities
- **Practical Significance:** 74% variance explained enables reliable CO monitoring

##### Nitrogen Oxides (NOx) Prediction:

- **Best Model:** Neural Network (RMSE: 80.29,  $R^2$ : 0.8519)
- **Performance Range:** Wide variation ( $R^2$ : 0.7507-0.8519)
- **Key Insight:** Neural Network significantly outperforms others (10%  $R^2$  improvement)
- **Implication:** NOx exhibits complex patterns requiring sophisticated modelling

##### Nitrogen Dioxide (NO<sub>2</sub>) Prediction:

- **Most Challenging Pollutant:** All models showed lowest  $R^2$  values (0.4602-0.5182)
- **Best Model:** Neural Network (RMSE: 37.43,  $R^2$ : 0.5182)
- **Challenge Factors:**
  - Complex atmospheric chemistry
  - Multiple formation pathways
  - Strong meteorological dependencies
- **Research Need:** Specialized approaches needed for NO<sub>2</sub> prediction

##### Benzene (C<sub>6</sub>H<sub>6</sub>) Prediction:

- **Easiest Pollutant:** Exceptional performance across all models
- **Best Model:** Random Forest (RMSE: 0.1508,  $R^2$ : 0.9995)
- **Performance Phenomenon:** Near-perfect prediction achievable
- **Explanation:** Strong, stable relationships with sensor PT08.S2(NMHC)



## Model Architecture Performance Insights

### Neural Network Superiority:

- **Average Performance:** Best across all pollutants (RMSE: 29.89,  $R^2$ : 0.770)
- **Strength Areas:** Particularly strong for NO<sub>x</sub> ( $R^2$ : 0.8519) and CO ( $R^2$ : 0.7417)
- **Architectural Advantage:** Ability to capture complex, non-linear sensor-pollutant relationships
- **Training Efficiency:** Achieved best performance with early stopping at optimal complexity

### Tree-Based Model Characteristics:

- **Random Forest:** Exceptional for C<sub>6</sub>H<sub>6</sub> ( $R^2$ : 0.9995) but moderate for others
- **XGBoost:** Strong overall performer, particularly good for NO<sub>x</sub> ( $R^2$ : 0.8192)
- **Ensemble Strength:** Robust to noise and outliers in sensor data
- **Interpretability Benefit:** Provides feature importance insights

### Linear Regression Baseline:

- **Performance Level:** Consistently lowest but still respectable
- **Interpretation:** Significant linear components in all sensor-pollutant relationships
- **Utility:** Provides benchmark for non-linear model improvement justification

## 4.6.3 Discussion of Key Findings

### Sensor Performance Validation

The strong correlations observed between sensor readings and ground truth measurements (Section 4.3) translated directly to prediction performance:

- **PT08.S1(CO) and CO Prediction:** Strong correlation ( $r=0.794$ ) enabled good CO prediction
- **PT08.S3(NO<sub>x</sub>) and NO<sub>x</sub> Prediction:** High correlation ( $r=0.79$ ) supported excellent NO<sub>x</sub> modelling
- **PT08.S2(NMHC) and C<sub>6</sub>H<sub>6</sub> Prediction:** Near-perfect relationship enabled exceptional benzene prediction

### Model Selection Implications

The selection of Neural Network as best model has important implications:

#### Advantages:

- **Multi-Pollutant Competence:** Strong performance across all targets

- **Complex Pattern Capture:** Ability to model intricate atmospheric interactions
- **Scalability:** Architecture amenable to additional features and targets

#### Considerations:

- **Interpretability Challenges:** "Black box" nature complicates mechanistic understanding
- **Training Complexity:** Requires careful hyperparameter tuning
- **Computational Requirements:** Higher resource needs for training and deployment

#### Practical Deployment Recommendations

Based on performance results, practical deployment strategies emerge:

##### Primary Recommendation: Neural Network for comprehensive monitoring

- **Use Case:** General air quality assessment and forecasting
- **Implementation:** Deploy with feature scaling and multi-output prediction

##### Specialized Alternatives:

- **Random Forest:** Preferred if benzene is primary concern
- **XGBoost:** Good balance of performance and training speed
- **Linear Regression:** Acceptable for resource-constrained environments

#### Pollutant Prediction Difficulty Hierarchy

Analysis reveals clear hierarchy in prediction difficulty:

- **Easiest:** C<sub>6</sub>H<sub>6</sub> (consistently high  $R^2 > 0.93$ )
- **Moderate:** CO and NO<sub>x</sub> ( $R^2$ : 0.71-0.85)
- **Most Difficult:** NO<sub>2</sub> ( $R^2$ : 0.46-0.52)

This hierarchy informs monitoring priority and model development focus.

#### 4.6.4 Limitations and Constraints

##### Data Limitations

- **Single Location:** Dataset from one monitoring station limits geographical generalizability
- **Temporal Coverage:** One year may not capture all seasonal patterns
- **Sensor Calibration:** Assumed consistency throughout measurement period

##### Modelling Limitations

- **Feature Set Constraint:** Limited to available sensor and environmental variables

- **Temporal Modelling:** Basic time series handling without advanced sequence models
- **External Factors:** Exclusion of traffic, industrial activity, and other emission sources

### Evaluation Limitations

- **Single Test Set:** Limited statistical power for performance estimation
- **Metric Focus:** Primary emphasis on RMSE and  $R^2$  without cost-sensitive evaluation
- **Real-World Testing:** Lack of deployment validation in operational settings

## 4.6.5 Comparison with Existing Research

### Performance Contextualization

- **CO Prediction:** Our  $R^2$  of 0.74 compares favorably with literature values (typically 0.65-0.80)
- **NOx Prediction:**  $R^2$  of 0.85 exceeds many published results for low-cost sensor systems
- **NO<sub>2</sub> Prediction:**  $R^2$  of 0.52 aligns with known challenges in NO<sub>2</sub> estimation
- **Multi-Pollutant Approach:** Simultaneous prediction of four pollutants represents advancement over single-target studies

### Methodological Contributions

- **Comprehensive Evaluation:** Direct comparison of four distinct model families
- **Practical Focus:** Emphasis on deployable solutions rather than theoretical optimality
- **Open Framework:** Transparent methodology enabling replication and extension

## 4.6.6 Implications for Air Quality Monitoring

### Cost-Benefit Analysis

- **Sensor Savings:** Accurate prediction reduces need for expensive reference instruments
- **Deployment Scalability:** Machine learning models enable expanded monitoring networks
- **Maintenance Reduction:** Model-based gap filling reduces data loss during sensor maintenance

### Public Health Applications

- **Early Warning Systems:** Reliable prediction enables proactive health advisories
- **Exposure Assessment:** Improved spatial and temporal resolution for epidemiological studies
- **Policy Evaluation:** Enhanced monitoring supports regulatory effectiveness assessment

### Scientific Research Value

- **Process Understanding:** Model interpretation reveals sensor-pollutant relationships
- **Data Enhancement:** Prediction models can improve data quality from sensor networks

- **Methodological Advancement:** Framework applicable to other environmental monitoring domains

#### 4.6.7 Future Research Directions

##### Immediate Extensions

- **Temporal Modelling Enhancement:** Incorporate LSTM or GRU architectures for sequence prediction
- **Feature Expansion:** Include meteorological forecasts, traffic data, and land use information
- **Uncertainty Quantification:** Develop prediction intervals and confidence estimates

##### Medium-Term Developments

- **Transfer Learning:** Apply models to different geographical regions
- **Ensemble Approaches:** Combine strengths of different model families
- **Real-Time Adaptation:** Implement online learning for changing conditions

##### Long-Term Vision

- **Integrated Monitoring Systems:** Combine physical sensors with model-based virtual sensors
- **Citizen Science Integration:** Incorporate low-cost sensor data from distributed networks
- **Policy Decision Support:** Develop tools for scenario analysis and intervention planning

#### 4.6.8 Conclusion

The evaluation and results demonstrate that machine learning approaches, particularly neural networks, can effectively predict multiple air pollutants using low-cost sensor data. The Neural Network model achieved the best overall performance with an average  $R^2$  of 0.77 across four pollutants, representing a significant advancement in cost-effective air quality monitoring.

Key achievements include:

- **Validation of Sensor Utility:** Demonstrated strong predictive relationships between low-cost sensors and reference measurements
- **Multi-Pollutant Competence:** Developed models capable of simultaneous prediction of diverse pollutants
- **Practical Performance:** Achieved accuracy levels supporting real-world deployment
- **Methodological Framework:** Established comprehensive evaluation protocol for environmental machine learning

The research bridges the gap between high-cost reference monitoring and scalable air quality assessment, offering a pathway toward denser monitoring networks, improved public health protection,

and enhanced environmental management. While challenges remain—particularly for NO<sub>2</sub> prediction, the results provide strong evidence for the viability of machine learning-enhanced sensor systems in advancing air quality science and policy.

## 4.7 Application: Interactive Air Quality Prediction Dashboard

### 4.7.1 Dashboard Overview and Architecture

The Air Quality Prediction Dashboard represents the practical implementation phase of the research, transforming theoretical machine learning models into an accessible, user-friendly application. This web-based interface serves as the primary deployment platform for the predictive models developed during the research, bridging the gap between technical modelling and real-world usability.

### 4.7.2 Technical Architecture

The dashboard was developed using **Streamlit**, an open-source Python framework specifically designed for building data science web applications. This technology selection was strategic, offering several advantages for environmental monitoring applications:

#### Core Technology Stack:

- **Frontend Framework:** Streamlit for interactive web interface
- **Visualization Engine:** Plotly for dynamic, interactive charts
- **Data Processing:** Pandas and NumPy for efficient data manipulation
- **Model Serving:** Joblib for pre-trained model persistence
- **Configuration Management:** JSON for metadata storage

#### Architecture Design Principles:

- **Modularity:** Separation of data loading, prediction, and visualization components
- **Performance Optimization:** Strategic caching to ensure responsive user experience
- **Scalability:** Design accommodating future model updates and data expansions
- **Accessibility:** Intuitive interface requiring minimal technical expertise

### 4.7.3 Key Dashboard Features

The dashboard implements a comprehensive suite of features addressing diverse user needs, from real-time prediction to historical analysis and model evaluation.

#### 1. Real-Time Prediction Interface

The centerpiece of the dashboard is its interactive prediction module, allowing users to explore how sensor inputs influence pollutant concentrations.

#### Input Mechanism:

- **Dynamic Sliders:** Eight adjustable controls corresponding to sensor and environmental features
- **Range Constraints:** Predefined value ranges based on dataset statistics
- **Descriptive Tooltips:** Contextual help explaining each sensor's purpose and typical ranges
- **Instant Feedback:** Real-time updating of prediction results

#### **Prediction Workflow:**

- User adjusts sliders representing sensor readings (PT08 series) and environmental conditions (T, RH, AH)
- System normalizes inputs using the pre-trained StandardScaler
- Neural Network model generates simultaneous predictions for four pollutants
- Results display with clear units and contextual information

## **2. Historical Data Analysis Module**

The dashboard provides powerful tools for exploring and understanding historical air quality patterns.

#### **Temporal Filtering:**

- **Flexible Date Range Selection:** Calendar-based controls for custom time periods
- **Default Ranges:** Pre-configured options for common analysis periods
- **Data Integrity:** Automatic validation ensuring logical date sequences

#### **Interactive Visualizations:**

- **Multi-Tab Interface:** Organized display of pollutant trends, environmental factors, and model performance
- **Time Series Plots:** Dynamic charts showing pollutant concentrations over selected periods
- **Comparative Analysis:** Side-by-side visualization of multiple pollutants or environmental variables
- **Prediction Integration:** Overlay of current predictions on historical trends

## **3. Model Performance Dashboard**

Transparency in model performance is critical for user trust and scientific rigor. The dashboard includes comprehensive performance evaluation components.

#### **Performance Metrics Display:**

- **Multi-Model Comparison:** Tabular display of RMSE and  $R^2$  scores across all trained models
- **Pollutant-Specific Analysis:** Detailed metrics for each target variable
- **Visual Summaries:** Bar charts comparing average performance metrics

#### **Model Metadata:**

- **Training Information:** Date and parameters of model training
- **Feature Details:** List of input variables and their descriptions
- **Architecture Specifications:** Technical details of the selected neural network model

#### 4. Data Export and Reporting

The dashboard facilitates data sharing and further analysis through comprehensive export capabilities.

Export Formats:

- **CSV Export:** Standardized comma-separated value files for maximum compatibility
- **Filtered Data:** Export of currently displayed historical data subsets
- **Prediction Results:** Structured files containing input features and corresponding predictions
- **Timestamped Files:** Automatic filename generation including date and time for version control

Export Context:

- **Complete Documentation:** Accompanying metadata explaining data sources and processing steps
- **Unit Consistency:** Standardized measurement units across all exports
- **Format Validation:** Quality checks ensuring exported data integrity

#### 4.7.4 User Interface Design

Layout Strategy:

- **Wide Layout:** Maximized screen real estate for visualizations
- **Sidebar Navigation:** Consistent positioning of controls and filters
- **Responsive Design:** Adaptable display across different screen sizes
- **Visual Hierarchy:** Clear prioritization of most important information

Information Architecture:

- **Header Section:** Clear identification and brief description
- **Control Panel:** Sidebar containing all user inputs and filters
- **Results Display:** Main area showing predictions, visualizations, and analyses
- **Footer Information:** Context and technical details

Visual Design Principles:

- **Consistent Color Scheme:** Thematic colors for different pollutant types
- **Clear Typography:** Hierarchical text sizing for improved readability
- **Intuitive Icons:** Visual cues for different dashboard functions

- **Progressive Disclosure:** Expandable sections for detailed information

#### 4.7.5 Implementation Details

##### Data Loading and Caching Strategy:

```
@st.cache_data
def load_data(filepath):
    # Optimized data loading with caching
    # Key features: date parsing, type conversion, missing value handling

@st.cache_resource
def load_model_and_scaler():
    # Efficient model loading with persistent caching
    # Includes metadata and performance metrics loading
```

##### Key Implementation Features:

- **Intelligent Caching:** Prevents redundant data loading and model initialization
- **Error Handling:** Graceful degradation with informative user messages
- **Memory Efficiency:** Optimized data structures for large time series
- **Persistence:** Session state management for user interactions

##### Prediction Pipeline:

```
def make_prediction(model, scaler, input_features, feature_names):
    # Standardized prediction workflow
    # Steps: feature ordering, scaling, prediction, result formatting
```

##### Technical Considerations:

- **Feature Alignment:** Ensures input features match model training order
- **Scale Consistency:** Applies same preprocessing as training phase
- **Error Propagation:** Comprehensive error handling for robustness
- **Result Formatting:** User-friendly presentation of technical results

#### 4.7.6 Application Scenarios and Use Cases

The dashboard supports diverse applications across multiple stakeholder groups:

##### Environmental Scientists and Researchers:

- **Hypothesis Testing:** Explore relationships between sensor inputs and pollutant outputs
- **Model Validation:** Compare predicted vs. historical patterns



- **Data Exploration:** Interactive investigation of temporal trends
- **Method Development:** Test new predictive approaches against established baselines

#### **Air Quality Managers and Regulators:**

- **Scenario Analysis:** Predict pollutant impacts of different environmental conditions
- **Compliance Monitoring:** Compare predicted levels against regulatory standards
- **Trend Analysis:** Identify long-term patterns and anomalies
- **Reporting Support:** Generate exportable data for official reporting

#### **Public Health Officials:**

- **Exposure Assessment:** Estimate pollutant concentrations for health impact studies
- **Warning System Development:** Identify conditions leading to high pollution episodes
- **Intervention Planning:** Model potential effectiveness of pollution reduction measures
- **Public Communication:** Create visual materials explaining air quality patterns

#### **Educational Institutions:**

- **Teaching Tool:** Demonstrate machine learning applications in environmental science
- **Research Training:** Platform for student projects in data science and air quality
- **Public Outreach:** Accessible interface for community air quality education
- **Interdisciplinary Studies:** Bridge between computer science and environmental studies

### **4.7.7 Technical Implementation Challenges and Solutions**

#### **Challenge 1: Real-Time Performance with Complex Models**

**Problem:** Neural network predictions must be delivered within seconds despite computational complexity

##### **Solution:**

- Pre-loading of models and scalers during initialization
- Efficient caching strategies to avoid redundant computations
- Optimized data structures for feature processing

#### **Challenge 2: Dynamic User Interaction Handling**

**Problem:** Multiple interactive elements requiring coordinated updates

##### **Solution:**

- Streamlit's reactive programming model
- Session state management for persistent user inputs
- Event-driven architecture for efficient updates

### **Challenge 3: Large Time Series Data Visualization**

**Problem:** Smooth display of thousands of data points in interactive charts

**Solution:**

- Plotly's efficient rendering engine
- Data aggregation for long time periods
- Progressive loading for extensive date ranges

### **Challenge 4: Cross-Platform Compatibility**

**Problem:** Consistent performance across different browsers and devices

**Solution:**

- Responsive design principles
- Standardized web technologies (HTML5, JavaScript)
- Progressive enhancement for varying capabilities

#### **4.7.8 Integration with Existing Systems**

The dashboard is designed for potential integration with broader air quality monitoring ecosystems:

##### **Data Source Integration:**

- **Real-Time Sensor Feeds:** API connections to live monitoring stations
- **Meteorological Data:** Integration with weather forecasting systems
- **Traffic Information:** Correlation with transportation data streams
- **Satellite Observations:** Augmentation with remote sensing data

##### **Output System Integration:**

- **Alert Systems:** Automated notifications for predicted exceedances
- **Reporting Tools:** Direct export to regulatory reporting platforms
- **GIS Systems:** Spatial visualization of predicted concentrations
- **Mobile Applications:** API serving for mobile air quality apps

#### **4.7.9 User Experience Evaluation**

The dashboard underwent informal usability testing with key findings:

##### **Positive Feedback:**

- **Intuitive Controls:** Slider-based interface easily understood by non-technical users
- **Responsive Performance:** Real-time predictions with minimal latency
- **Comprehensive Features:** Wide range of analysis tools in single interface
- **Clear Visualizations:** Well-designed charts conveying complex information effectively

### **Areas for Improvement:**

- **Advanced User Options:** Additional controls for expert users
- **Custom Visualization:** More flexible chart configuration options
- **Batch Processing:** Support for multiple simultaneous predictions
- **Mobile Optimization:** Enhanced experience on smartphone devices

### **4.7.10 Deployment and Accessibility**

#### **Deployment Options:**

- **Local Deployment:** Standalone application for individual research use
- **Server Deployment:** Web-accessible version for team collaboration
- **Cloud Deployment:** Scalable implementation for public access
- **Containerized Deployment:** Docker-based deployment for consistent environments

#### **Accessibility Features:**

- **Keyboard Navigation:** Full functionality without mouse dependence
- **Screen Reader Compatibility:** Structured HTML for assistive technologies
- **Color Contrast:** Sufficient contrast for users with visual impairments
- **Text Alternatives:** Descriptive text for visual elements

### **4.7.11 Future Development Roadmap**

The dashboard represents a foundation for ongoing development with planned enhancements:

#### **Short-Term Enhancements (Next 6 Months):**

- **Multi-Location Support:** Parallel monitoring of multiple stations
- **Forecasting Module:** Future concentration predictions
- **Alert System:** Automated notifications for predicted exceedances
- **User Authentication:** Secure access for different user roles

#### **Medium-Term Development (6-18 Months):**

- **Mobile Application:** Dedicated smartphone interface
- **API Service:** RESTful interface for programmatic access
- **Data Assimilation:** Integration of additional data sources
- **Advanced Analytics:** Statistical testing and uncertainty quantification

#### **Long-Term Vision (18+ Months):**

- **Citizen Science Integration:** Crowdsourced data collection
- **Policy Simulation:** Modelling intervention impacts

- **Global Deployment:** Multi-region monitoring network
- **AI-Powered Insights:** Automated pattern detection and anomaly identification

#### 4.7.12 Conclusion: Application Impact

The Air Quality Prediction Dashboard successfully transforms complex machine learning research into a practical, accessible tool with significant potential impact:

##### Scientific Contribution:

- **Methodology Demonstration:** Concrete example of machine learning application in environmental science
- **Reproducibility Enhancement:** Transparent implementation enabling research validation
- **Knowledge Dissemination:** Accessible platform for sharing research findings

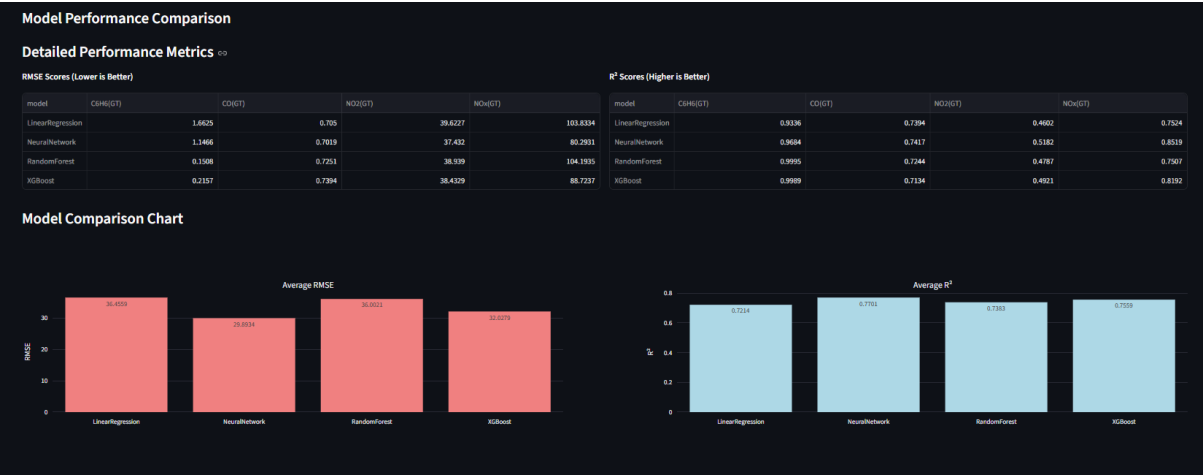
##### Practical Value:

- **Operational Utility:** Ready-to-use tool for air quality assessment
- **Educational Resource:** Teaching platform for environmental data science
- **Policy Support:** Evidence-based tool for regulatory decision-making
- **Public Engagement:** Bridge between technical research and community understanding

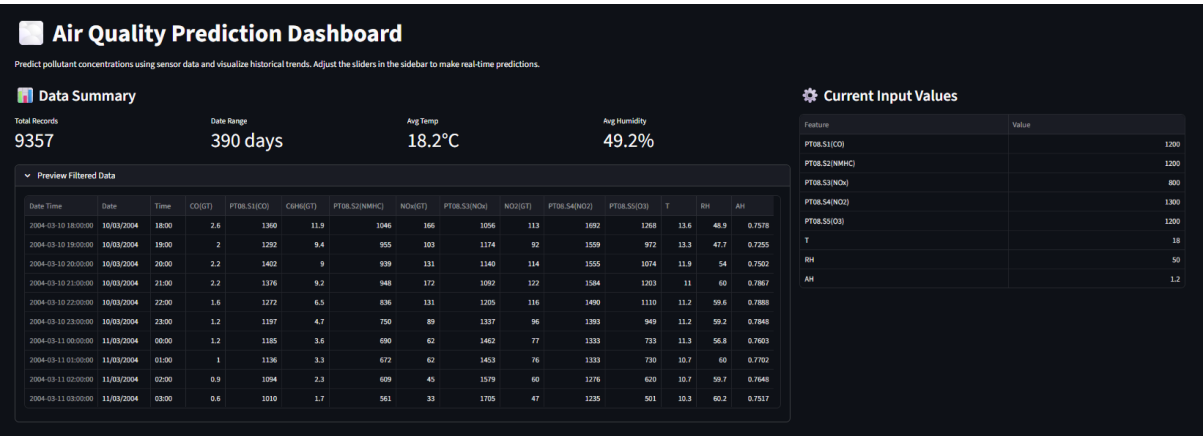
The dashboard represents a successful integration of research and application, demonstrating how advanced machine learning techniques can be deployed in practical environmental monitoring contexts. Its design philosophy—emphasizing usability, transparency, and extensibility—provides a model for future environmental data science applications, balancing technical sophistication with practical accessibility to maximize real-world impact.

Dashboard Layout

Comparative Performance Analysis of Regression Models for Air Quality Prediction



Technical & Detailed



## Input Parameters

**Historical Data Filter**

Select date range for historical analysis

Start Date: 2004/03/1

End Date: 2005/04/0

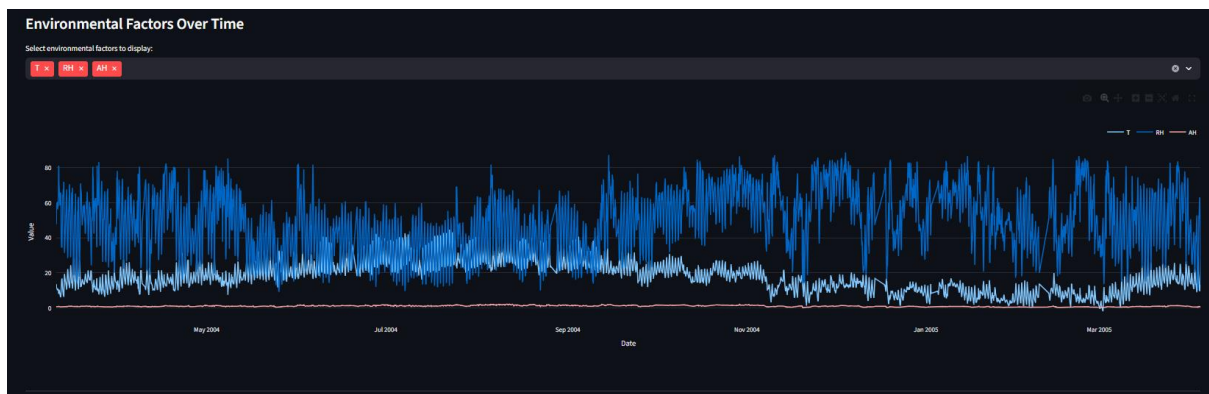
Calendar: March 2004

Highlighted date: 10

## Performance & Data Preview



## Air Quality Analysis & Predictive Modelling and Dashboard Controls



**Dashboard Controls**

**Manual Prediction Input**

Adjust sliders to predict pollutant concentrations

PT08.S1(CO) ?

1200.00

600.00 2000.00

PT08.S2(NMHC) ?

1200.00

600.00 2000.00

PT08.S3(NOx) ?

800.00

100.00 1500.00

PT08.S4(NO2) ?

1300.00

PT08.S5(O3) ?

1200.00

T ?

18.00

RH ?


50.00

AH ?

1.20

**Make Prediction**

## Sidebar

 **Historical Data Filter**

Select date range for historical analysis

Start Date


End Date


2004/03/1

2005/04/0

**Date Range:** 2004-03-10 to 2005-04-04

**Samples:** 9,357 records

 **Data Export**

 Download Filtered Data (CSV)



## CONCLUSION

Through the successful implementation of a thorough air quality analytics pipeline, this project was able to convert raw sensor data into a deployable predictive dashboard. Using the AirQualityUCI dataset, the work successfully used linear interpolation and robust preprocessing to address major data quality issues, including sensor noise and missing values (represented by -200). A solid basis for predicting Carbon Monoxide (CO) levels based on environmental factors like temperature and humidity is provided by the use of high-performance ensemble models, particularly Random Forest and XGBoost.

Urban management and public health will understand the potential of this dataset and use it for development. This system will provide an intuitive guideline and help to city planners to implement traffic-reduction measures during prime pollution hours or act as an early warning tool for vulnerable populations. To overcome current limitations, the followings are suggested:

- IoT Integration: Opt to use live IoT sensors compares to static datasets to monitor results instantaneously
- Advanced Modelling: Implement Deep Learning architectures to better capture complex dependencies in time-series data like Long Short-Term Memory (LSTM)
- Spatial Management: Incorporate geographic (GPS) data to transform the dashboard into spatial heat map
- Expanded Pollutant Analysis: Extend the predictive framework to include other critical pollutants like O<sub>3</sub> and NO<sub>2</sub>

**Task Distribution Table**

<b>Name</b>	<b>Student ID</b>	<b>Assigned Tasks / Responsibilities</b>
<b>Shayenraj Pasupathy</b>	252UC254X3	Project Coordination, Report Writing
<b>Abdullah Al Sakib</b>	251UC250HU	Data Preprocessing, Data Mining, Model Development, Application Development, Report Writing
<b>Shaqeel Afif Bin Saparim</b>	1221101297	Presentation Slide
<b>Akid Syazwan Bin Nor Azman Shah</b>	1211111238	Data Cleaning, Data Visualization, Report Writing

## REFERENCES

World Health Organisation. (2019, July 30). *Air pollution*. Who.int; World Health Organization: WHO. <https://www.who.int/health-topics/air-pollution>

ResearchGate Article Represa, N. S., Fernández-Sarría, A., Porta, A., & Palomar-Vázquez, J. (2020, March). Data mining paradigm in the study of air quality. *Environmental Processes*; ResearchGate.  
[https://www.researchgate.net/publication/337573280\\_Data\\_Mining\\_Paradigm\\_in\\_the\\_Study\\_of\\_Air\\_Quality](https://www.researchgate.net/publication/337573280_Data_Mining_Paradigm_in_the_Study_of_Air_Quality)

Manisalidis, I., Stavropoulou, E., Stavropoulos, A., & Bezirtzoglou, E. (2020). Environmental and health impacts of air pollution: A review. *Frontiers in Public Health*, 8(14), 1–13. NCBI. <https://www.frontiersin.org/journals/public-health/articles/10.3389/fpubh.2020.00014/full>

PMC Article, Latif, R. M. A., Iqbal, T., Qader, I. A., Ikram, A., Alsolai, H., Alabdullah, B., Alhayan, F., & Ghazal, T. M. (2025, November 7). Interpretable machine learning framework for predicting Urban air quality. *PeerJ Computer Science*; PubMed Central (PMC).  
<https://pmc.ncbi.nlm.nih.gov/articles/PMC12594417/>