

Investigating Air Quality Factors and Their Impact on Human Health

Group 28 – DSE 501 (Spring 2025)

Mohammad Hamza Choudhry • Kushal Trivedi

Kaviya Venkateshbabu • Srujana Rachamalla

Instructor: Prof. Rong Pan | Arizona State University

Agenda

- Introduction and Motivation
- Dataset Overview
- Data Cleaning and Preprocessing
- Exploratory Data Analysis
- Hypothesis Testing (H1–H6)
- Key Findings
- Public Health Implications and Policy Recommendations
- Future Work

Introduction

Air pollution remains a major global health challenge.

Urban populations are especially at risk due to:

- Heavy traffic and industrial emissions
- Seasonal factors like winter heating
- Weather conditions that trap pollutants near the ground

Our project focuses on:

- Analyzing when and why pollution peaks occur
- Evaluating the reliability of sensor-based air monitoring
- Connecting patterns in pollution to human health concerns

Dataset Overview

We used the UCI Air Quality dataset, which contains:

- 9,358 hourly records from an urban location in Italy
- Ground-truth pollutant values: CO, NO_x, NO₂, Benzene
- Sensor readings: 5 metal oxide sensors (PT08.S1–S5)
- Weather data: Temperature, Relative Humidity, Absolute Humidity
- Time features: Date and Time (used to extract hour, weekday, month)

This dataset allowed us to explore pollutant behavior, sensor accuracy, and environmental effects over time.

Air Quality Dataset: Attribute Definitions and Descriptions

The dataset includes hourly measurements from an air quality monitoring station. Each column represents a pollutant, sensor reading, or environmental variable.

Column Name	Description
Date	Date of measurement (DD/MM/YYYY)
Time	Time of measurement (HH.MM.SS)
CO(GT)	True CO concentration in mg/m^3
PT08.S1(CO)	Sensor 1 output in response to CO
NMHC(GT)	Non-Methane Hydrocarbons in $\mu\text{g/m}^3$ (many values missing)
C6H6(GT)	Benzene concentration in $\mu\text{g/m}^3$
PT08.S2(NMHC)	Sensor 2 output targeting NMHC
NOx(GT)	True NOx concentration in ppm

Air Quality Dataset: Attribute Definitions and Descriptions

Column Name	Description
PT08.S3(NOx)	Sensor 3 output targeting NOx
NO2(GT)	True NO ₂ concentration in µg/m ³
PT08.S4(NO2)	Sensor 4 output targeting NO ₂
PT08.S5(O3)	Sensor 5 output targeting O ₃
T	Ambient temperature in °C
RH	Relative humidity (%)
AH	Absolute humidity in g/m ³
Hour / Month	Extracted from timestamp for temporal trend analysis

This structured dataset enabled detailed analysis of pollutant behavior, sensor reliability, and temporal trends.

Data Cleaning & Preprocessing

To ensure data quality and consistency, we performed the following steps:

- Replaced all invalid values (e.g., -200) with missing value markers
- Dropped incomplete rows to retain only valid hourly records
- Combined Date and Time columns to form a single **Datetime** column
- Extracted new features:
 - **Hour of day**
 - **Day of week**
 - **Month**
- Binned temperature and humidity values to analyze joint effects on pollutants

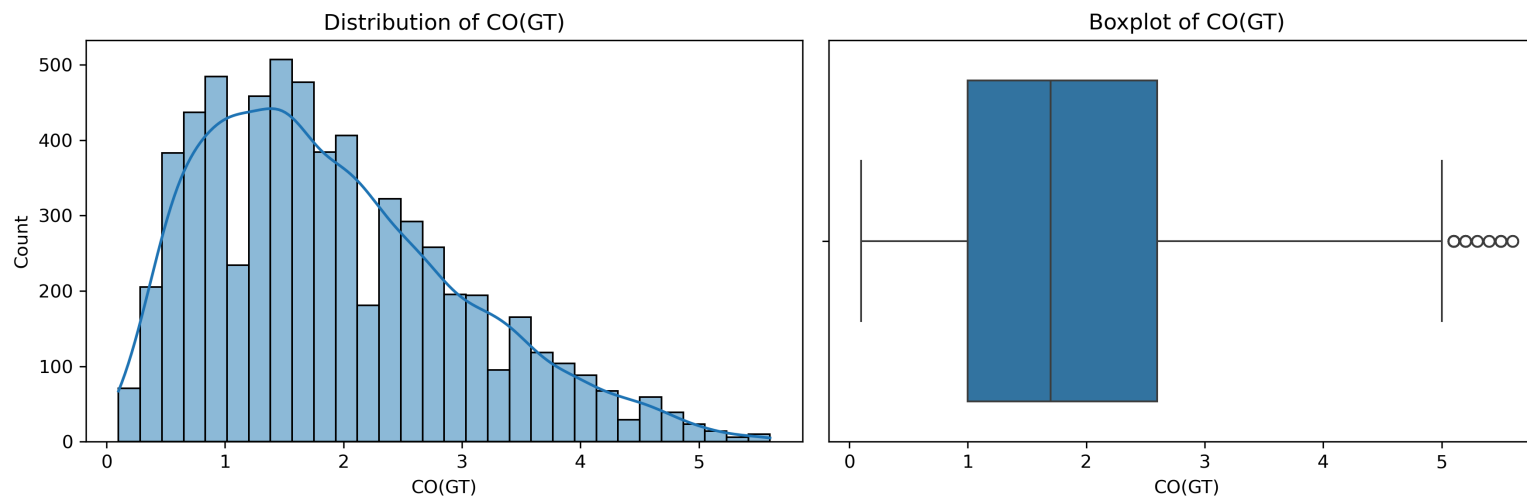
These steps prepared the dataset for time-based analysis and hypothesis testing.

Pollutant Distributions

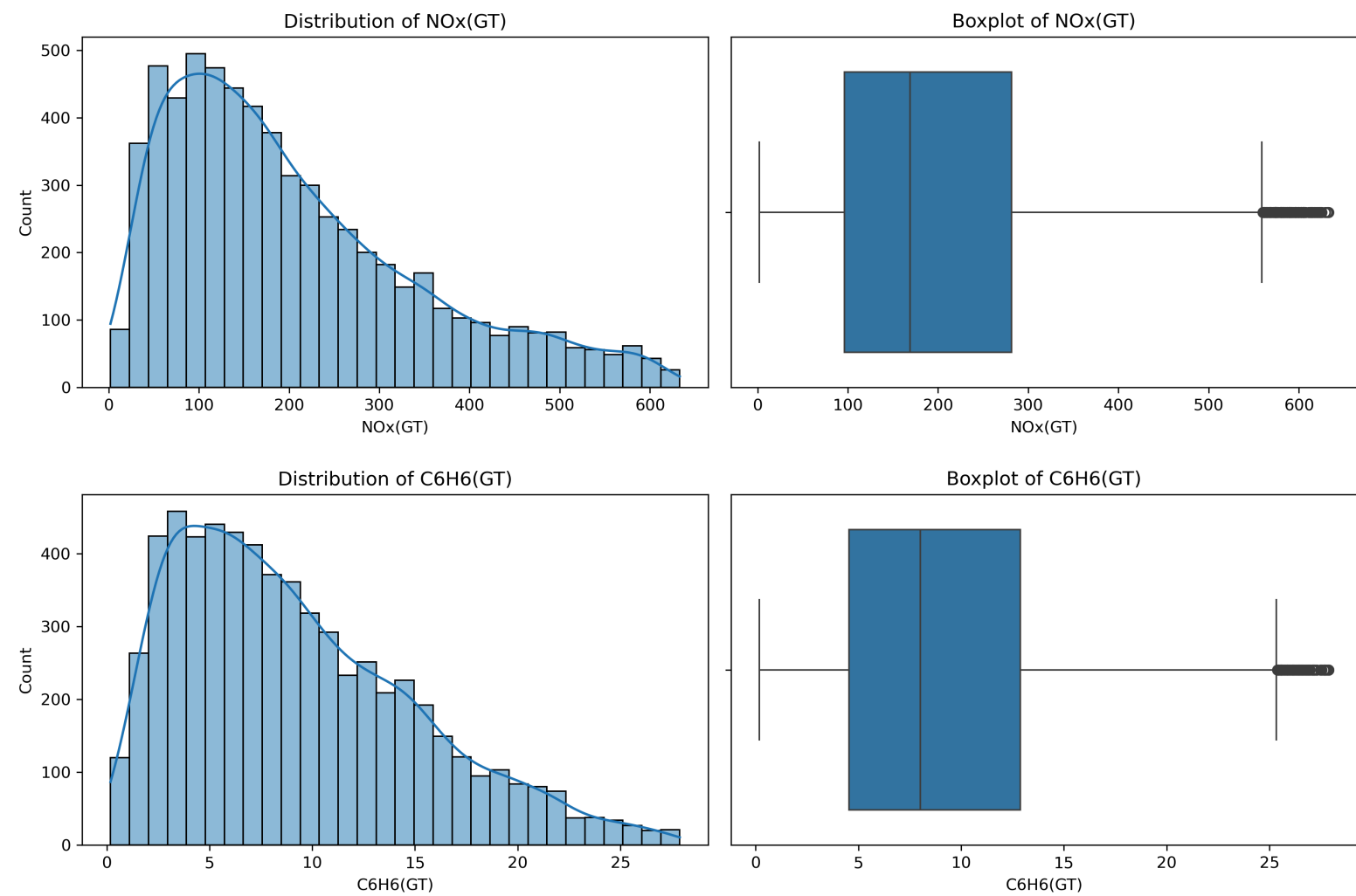
We examined the distribution of each key pollutant:

- Most pollutants showed **right-skewed patterns**
- A majority of values were low to moderate
- Some **high outliers** reflect traffic or industrial spikes

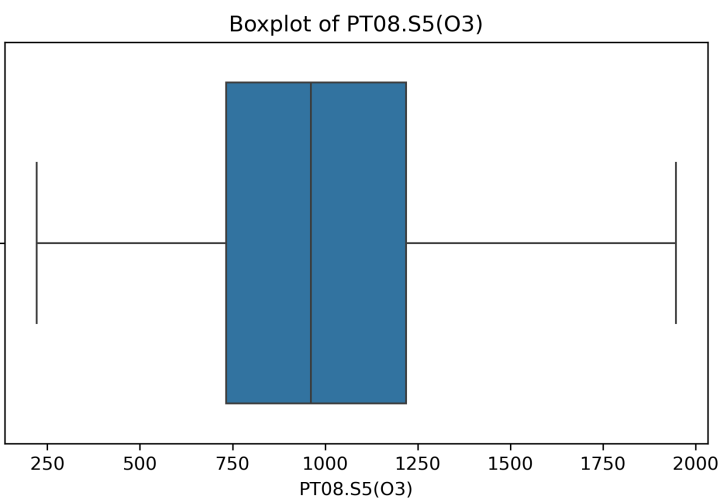
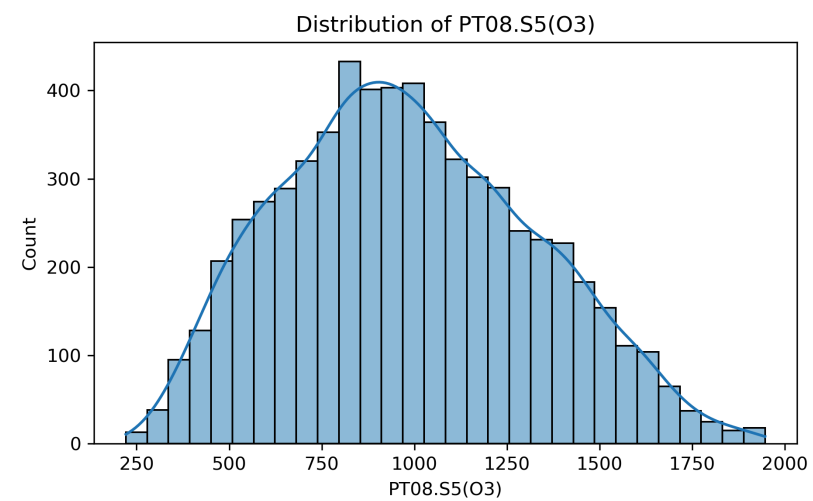
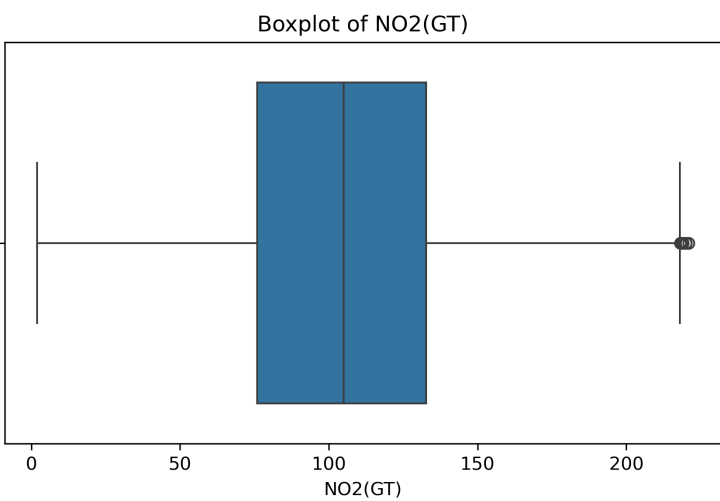
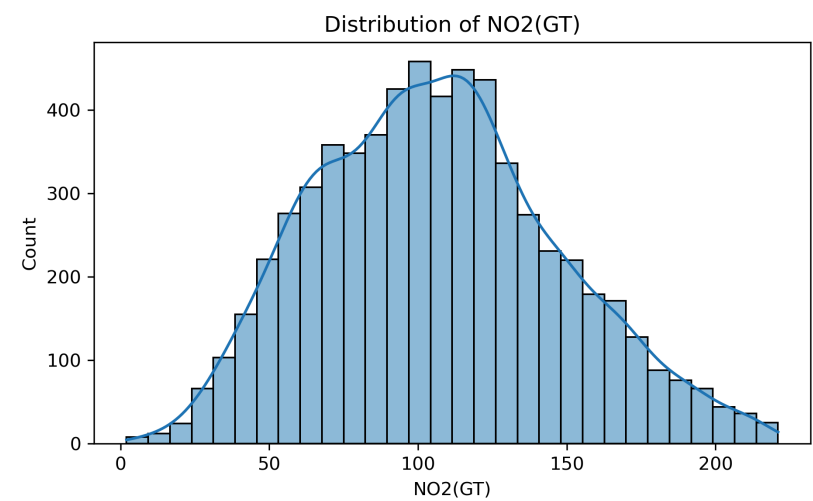
Visuals:



Pollutant Distributions



Pollutant Distributions



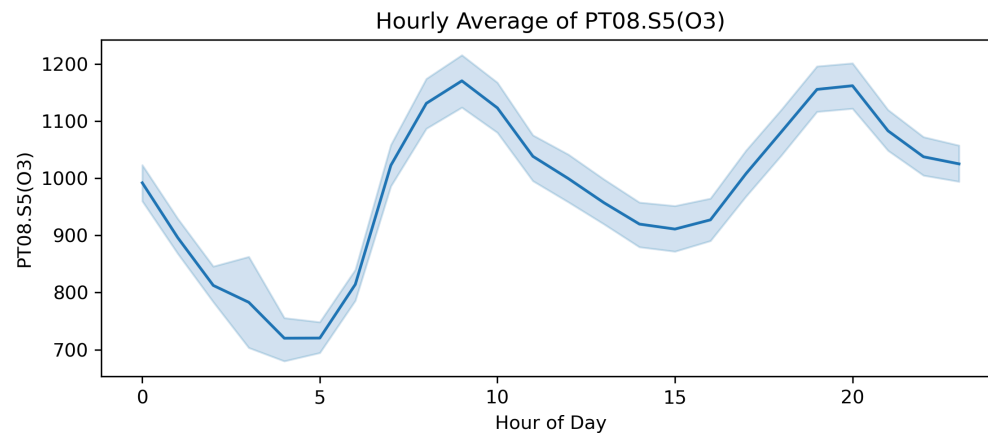
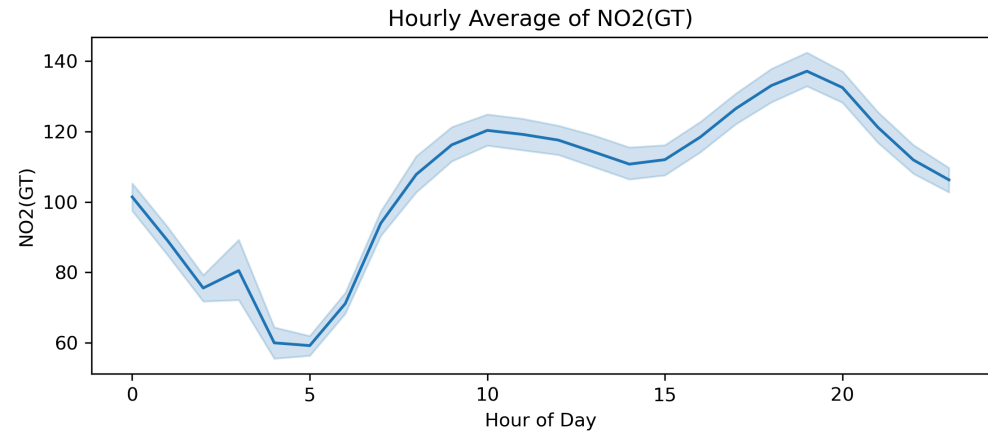
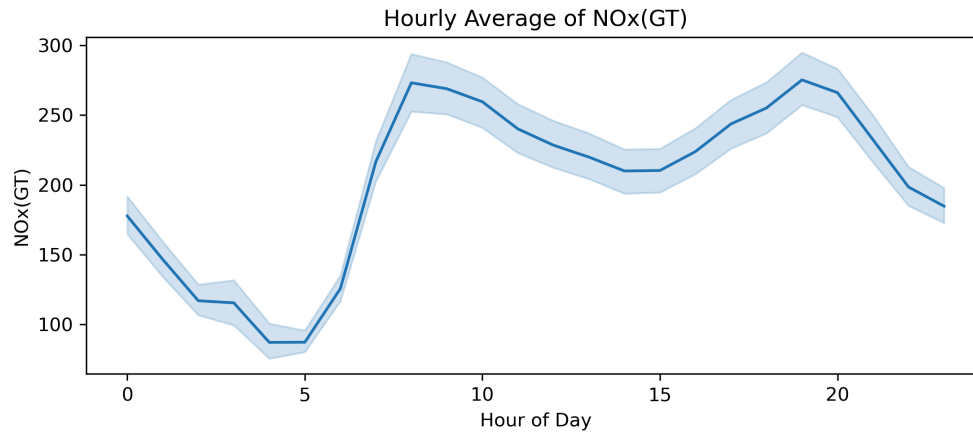
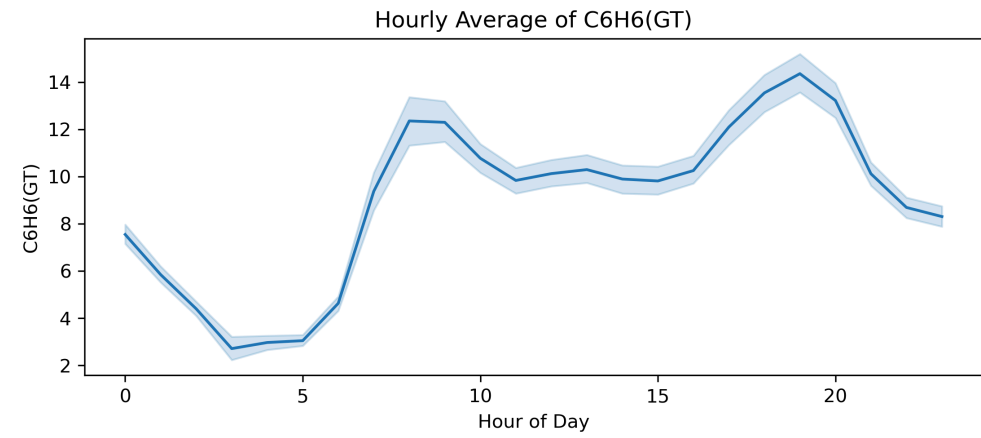
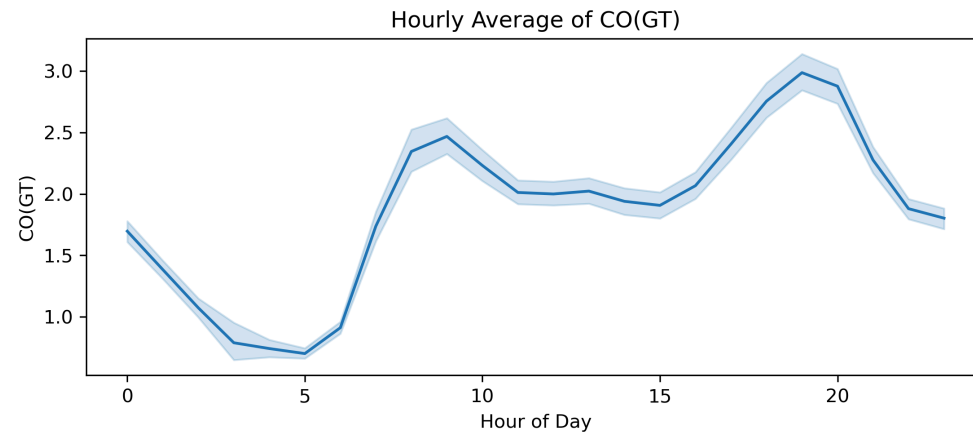
Daily and Weekly Trends

We observed clear temporal patterns in pollutant levels:

- **Rush hour peaks** around 8 AM and 6 PM for CO, NO_x and C₆H₆
- **Weekday levels** are consistently higher than weekends
- **Ozone (O₃)** tends to rise during mid-day sunlight hours

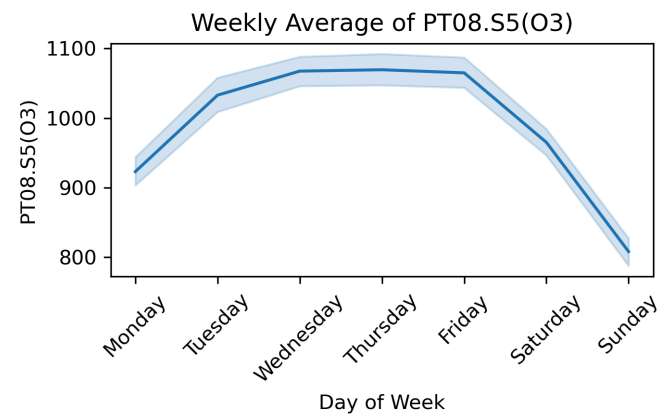
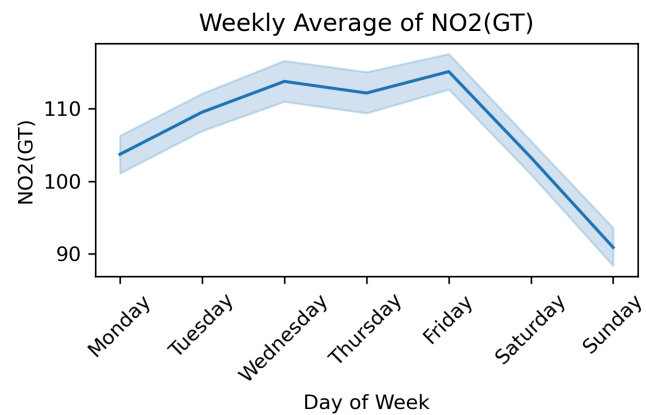
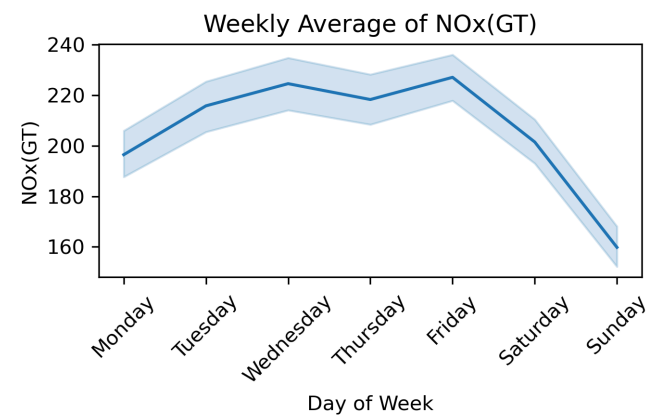
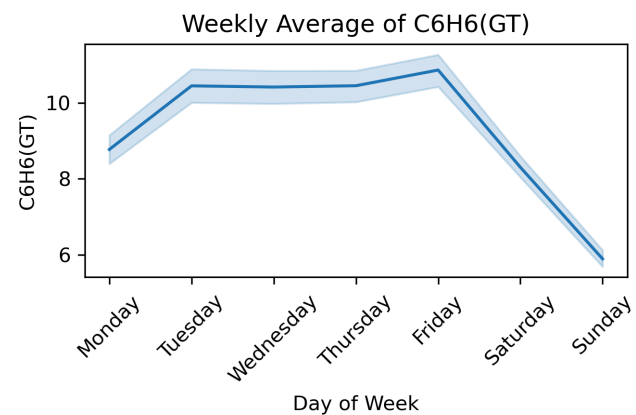
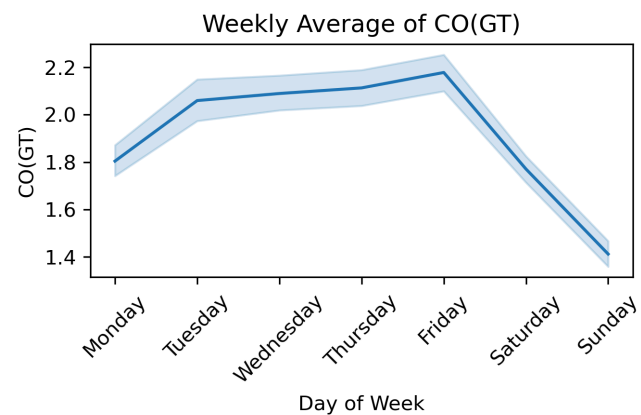
Hourly Trends: Pollution Peaks

- CO, NO_x, NO₂, and Benzene (C₆H₆) show clear peaks around 8 AM and 6 PM, aligning with rush hours
- O₃ levels rise steadily in the morning and **peak mid-day**, consistent with **sunlight-driven reactions**



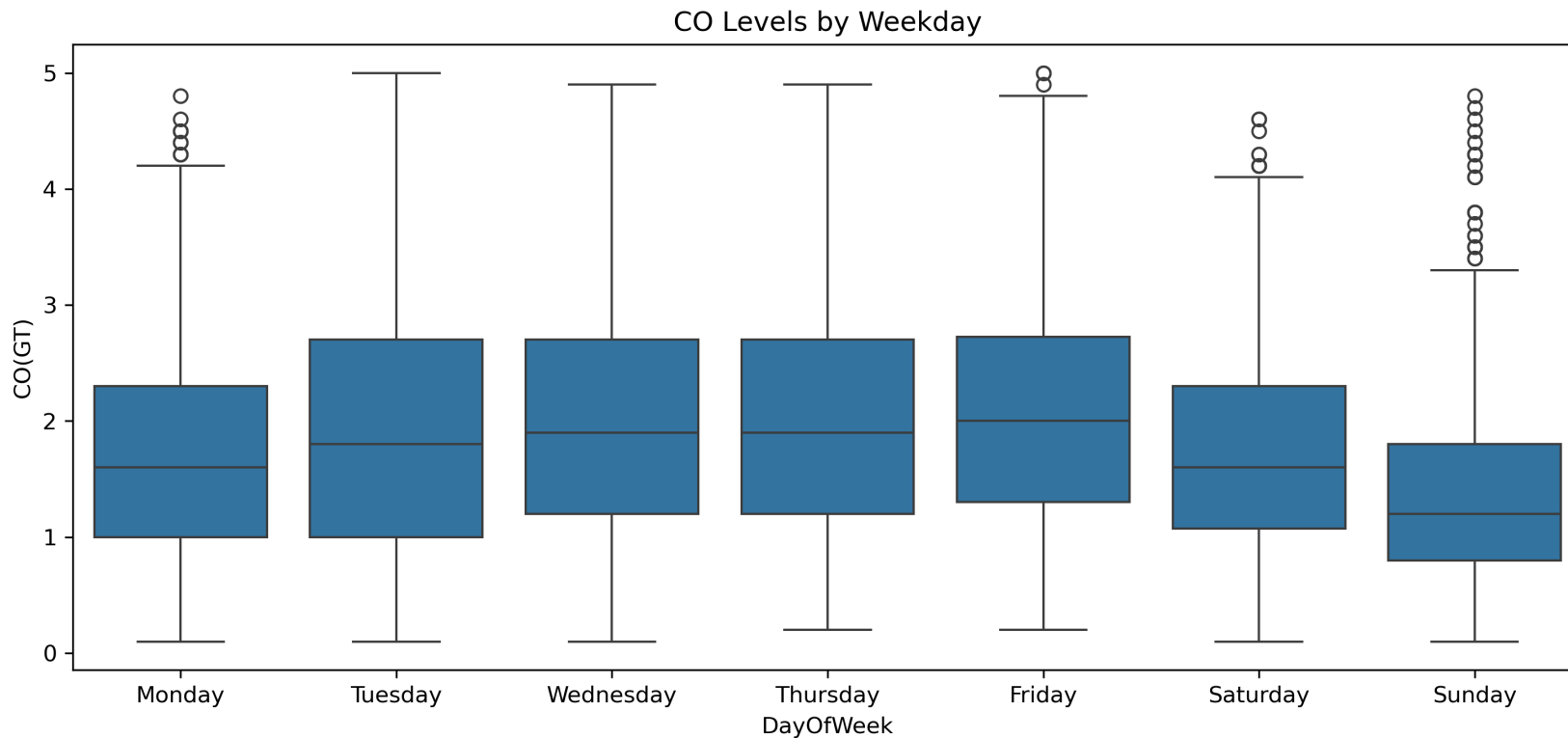
Weekly Trends: Weekdays vs. Weekends

- Pollutant levels are consistently **higher on weekdays**, especially Tuesday to Friday
- Sharp drop on **Saturdays and Sundays**, likely due to reduced traffic and industrial output



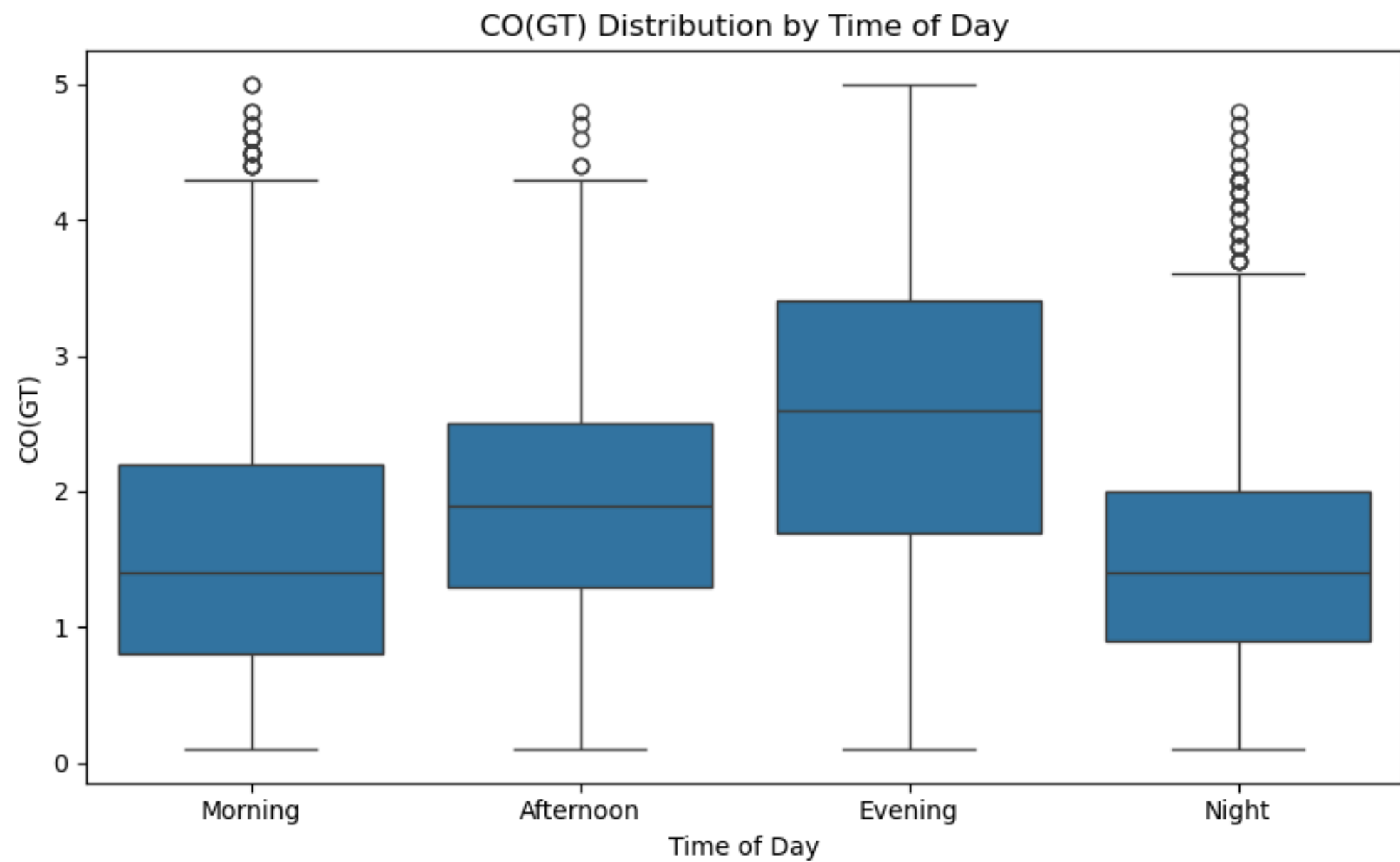
Weekly CO Variation: Boxplot View

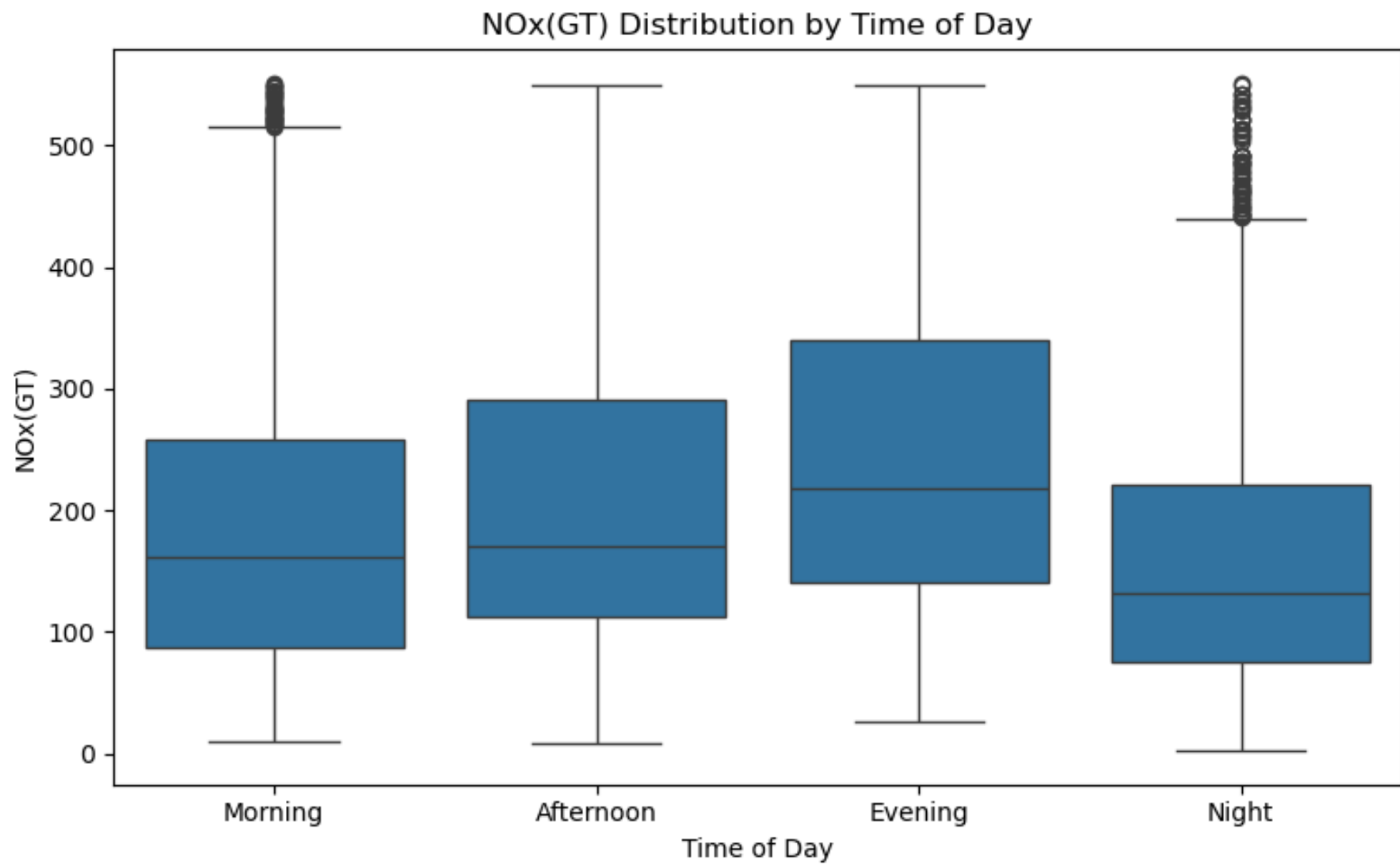
- CO concentrations are higher midweek and lower on weekends
- **Sunday** shows the **cleanest air** overall
- The boxplot reveals variability and confirms traffic-related pollution trends



Time of Day Patterns: CO and NO_x

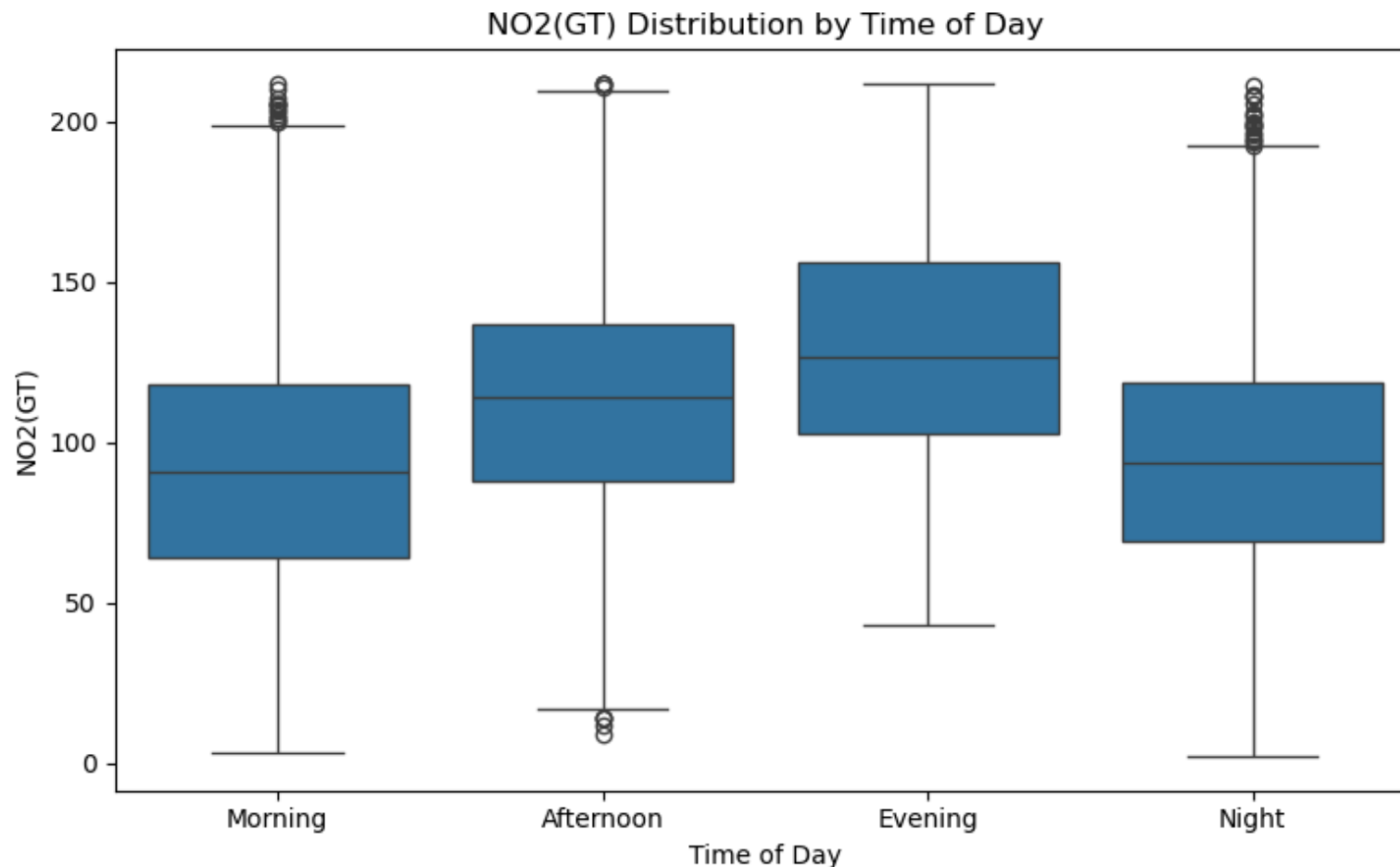
- Both **CO** and **NO_x** concentrations peak during the **evening** period.
- This reflects typical urban patterns tied to **evening rush hour** and reduced dispersion after sunset.
- CO levels are slightly higher in the afternoon than in the morning, but **evening is the dominant peak**.



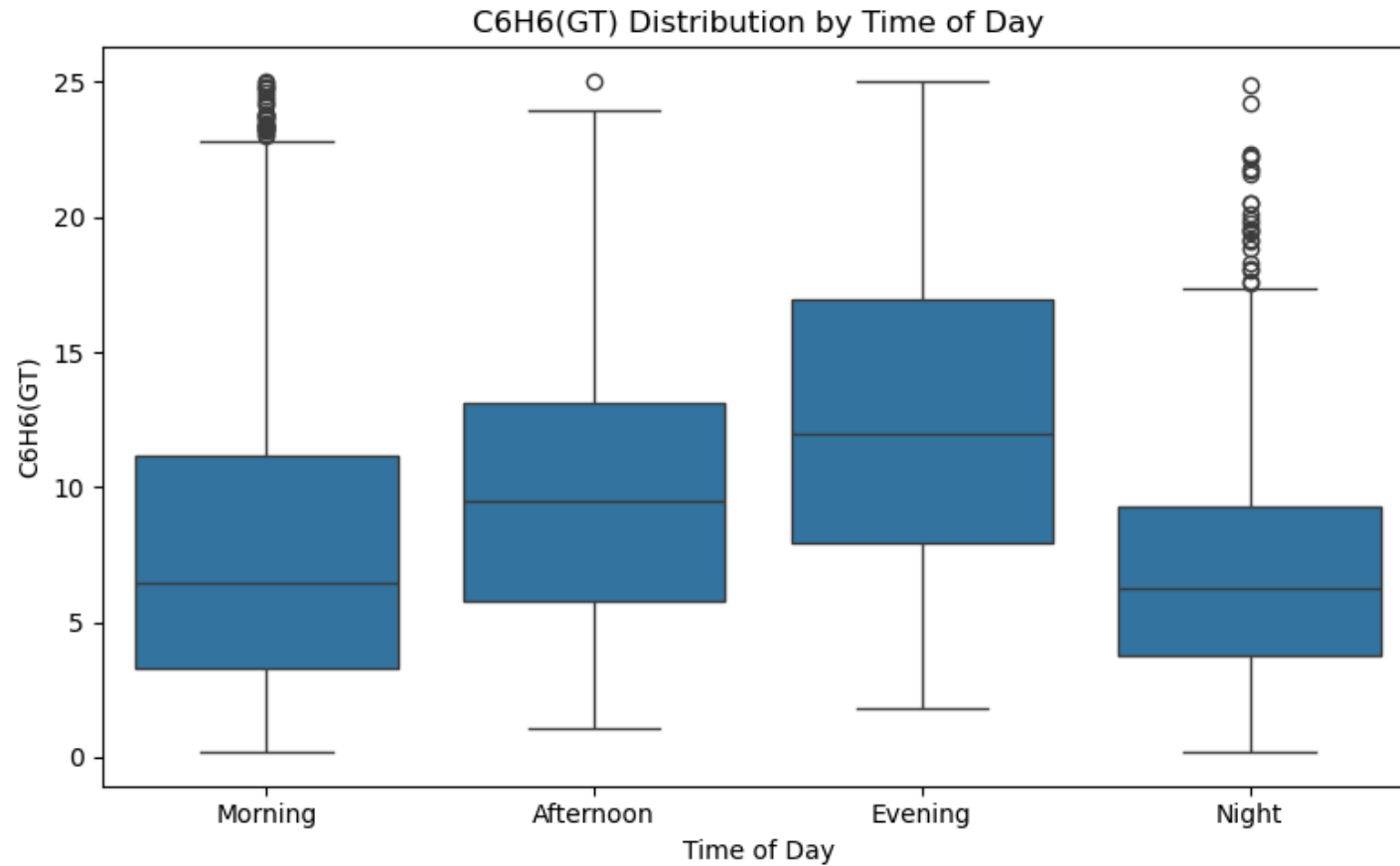


Time of Day Patterns: NO₂

- NO₂ concentrations increase throughout the day, reaching a **clear peak in the evening**.
- Follows similar dynamics to NO_x, suggesting shared traffic-based sources.

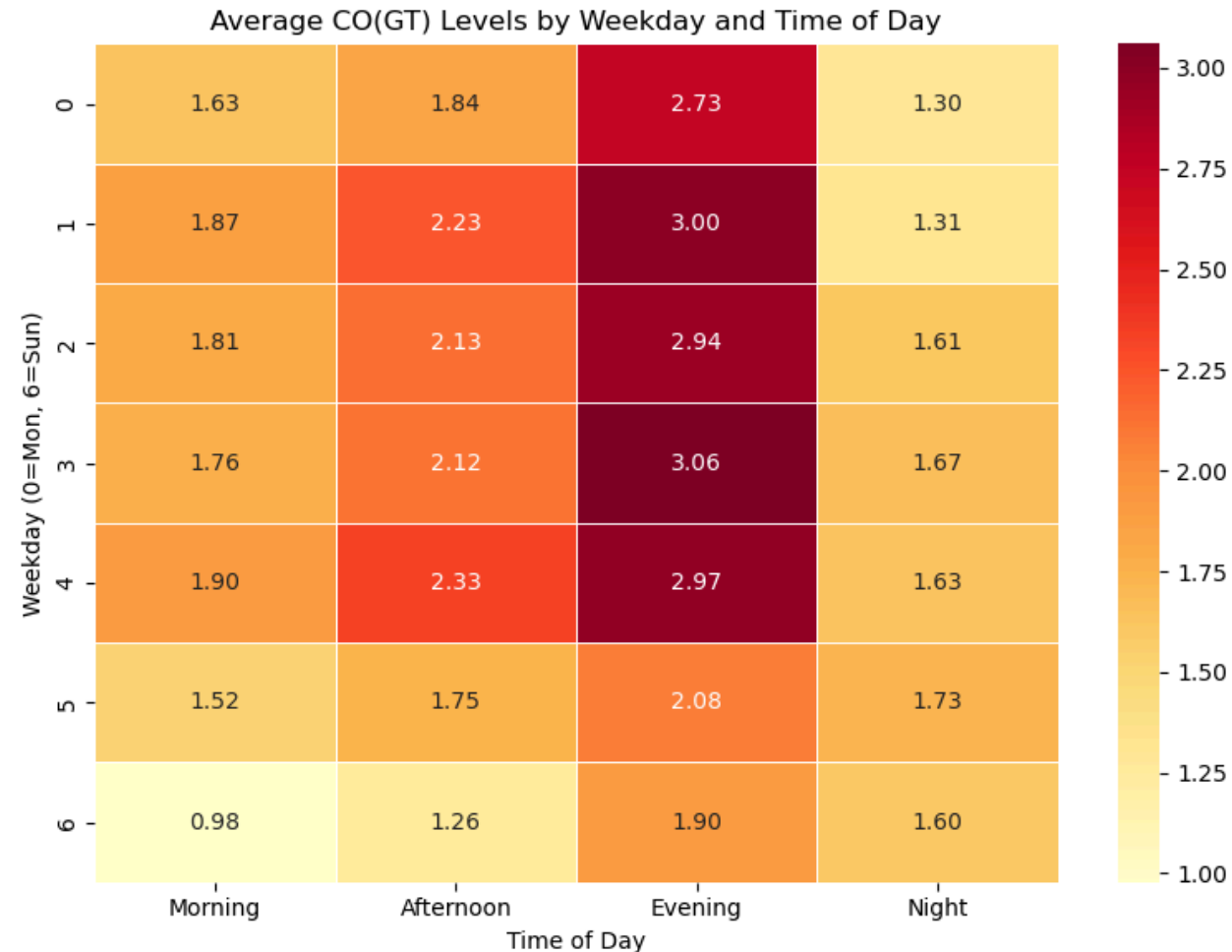


- Benzene (C_6H_6) also peaks in the **evening**, likely due to workday industrial and traffic emissions.



Cross-Time Heatmap: Weekday × Hour

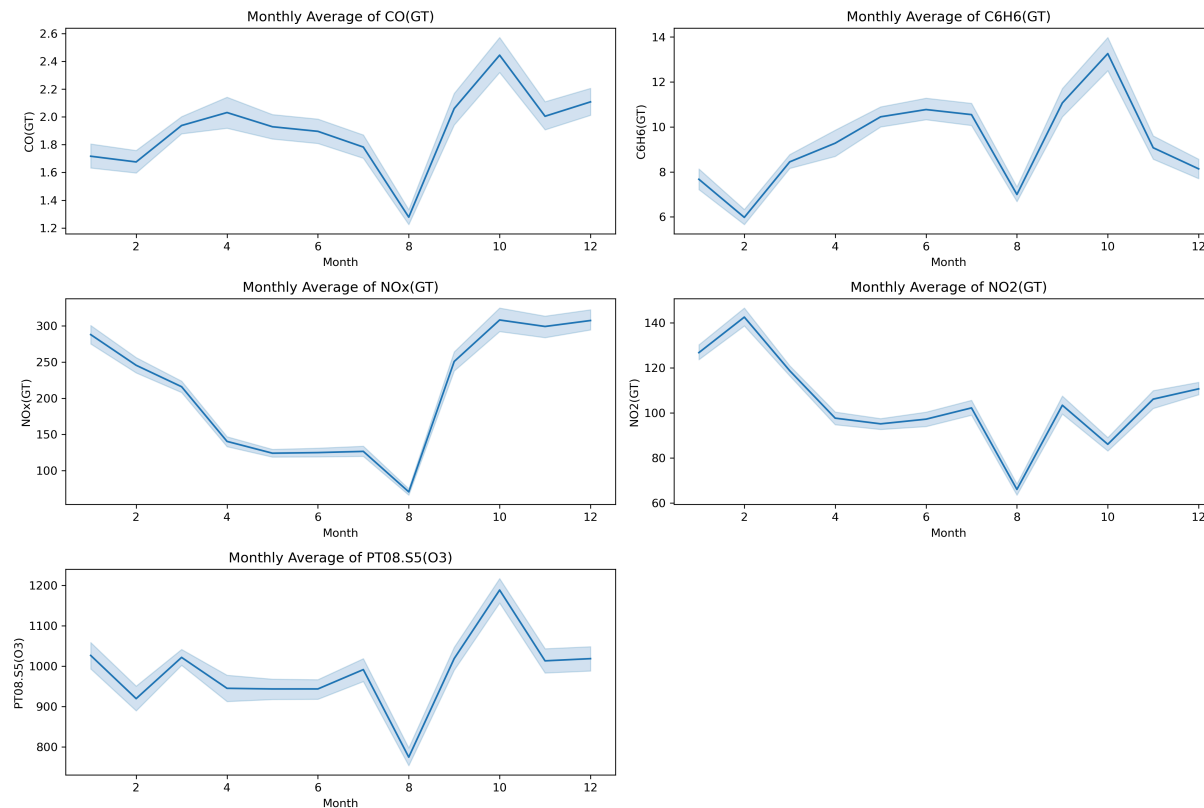
Highest CO concentrations are observed on **weekday evenings** — especially Tuesday to Thursday.



Seasonal Pollution Trends

Pollutants like CO and NO_x are significantly higher in **winter months**.

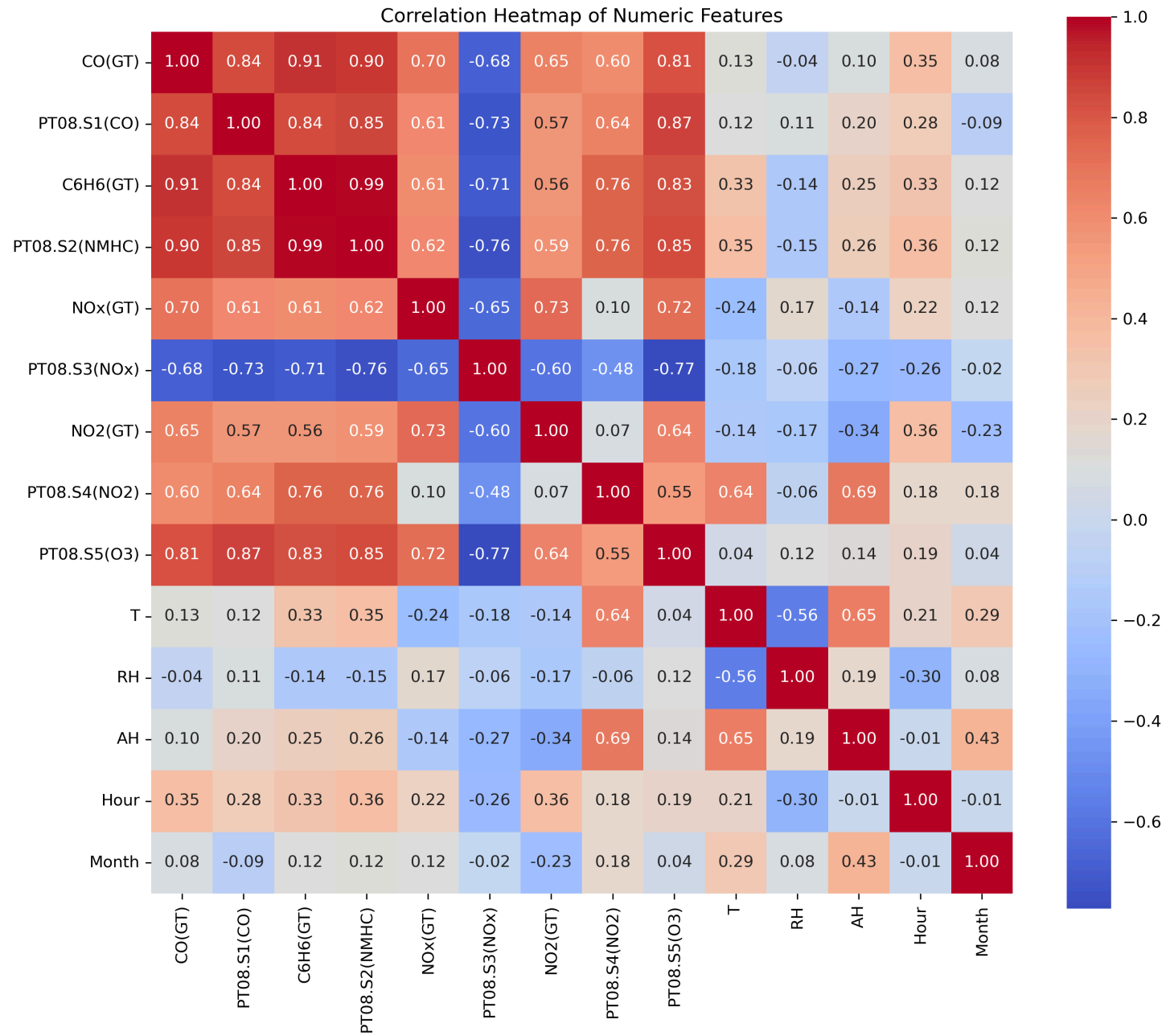
- August shows a marked decline in all pollutants, possibly from reduced industrial activity and better atmospheric dispersion. These trends highlight seasonal variation influenced by weather and human activity.



Correlation Analysis: Sensor Validity and Environmental Influence

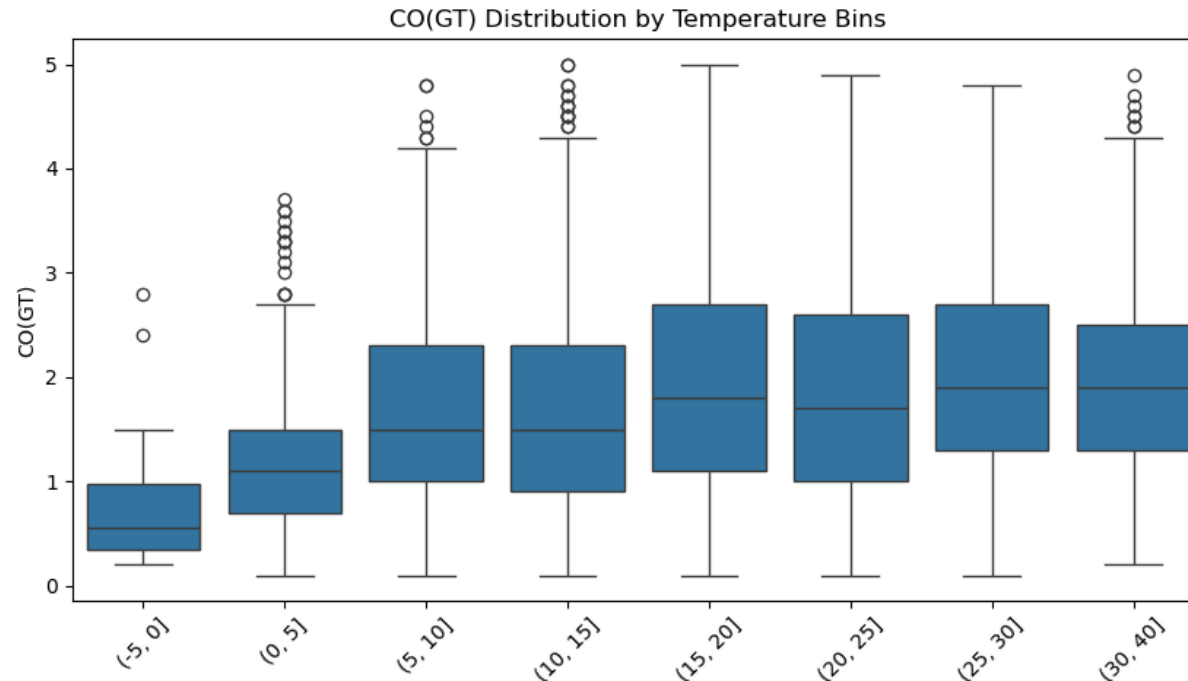
We analyzed how well each sensor correlates with its corresponding ground-truth pollutant. Key insights:

- **CO Sensor (PT08.S1):** Shows a **strong positive correlation** with CO(GT) ($r = 0.84$).
→ Indicates reliable performance and consistent tracking of CO levels.
- **NO_x Sensor (PT08.S3):** Displays a **strong negative correlation** with NO_x(GT) ($r = -0.65$).
→ This inverse relationship suggests a **serious calibration issue**, making the sensor unreliable for real-world measurements.
- **O₃ Sensor (PT08.S5):** Shows weak correlation with temperature ($r = 0.04$) and humidity ($r = 0.12$). Suggests that O₃ behavior in this dataset is influenced by additional environmental or chemical factors not captured here.

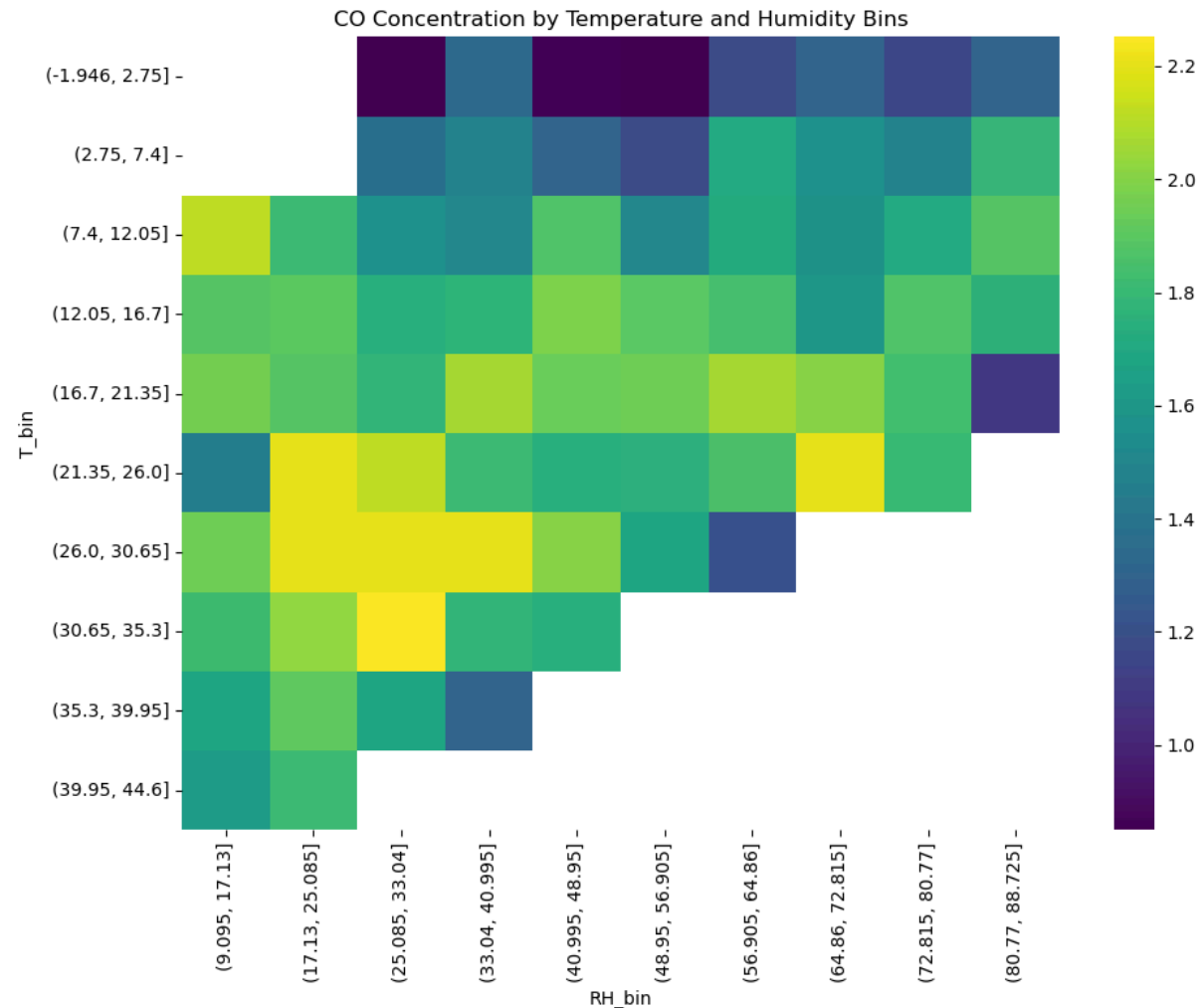


Temperature and Humidity Effects on CO

- CO levels increase with temperature until about 25°C, then plateau.
- This supports the idea that **cool to mild temperatures trap pollutants** closer to the ground.
- Outliers in these bins may represent **morning traffic** or **evening inversions**.



- CO concentration is **not linearly dependent** on temperature or humidity.
- Highest values cluster around **20–30°C** and **40–70% RH**
- This suggests pollutants **accumulate under moderate weather conditions** — possibly due to human activity and limited dispersion.



Hypothesis Testing

Our statistical investigation was guided by six research hypotheses.

Each hypothesis aimed to explore a distinct environmental or technical factor based on air quality literature and sensor evaluation objectives.

- **H1: Weekday Rush Hour Impact**

CO and NO_x levels are significantly higher during weekday rush hours compared to non-peak hours, reflecting traffic congestion patterns.

- **H2: Seasonal Variation in Pollutants**

Pollution levels for CO and NO_x are elevated in winter months due to low atmospheric dispersion and increased combustion-related heating.

- **H3: Influence of Temperature and Humidity**

Meteorological variables modulate pollutant behavior:

- O₃ increases with temperature
- CO and NO_x concentrations decrease with relative humidity

- **H4: Weekly Trends in Benzene (C₆H₆)**

Benzene concentrations vary significantly by day of the week, possibly due to variations in weekday industrial or traffic-related activity.

- **H5: Sensor–Ground Truth Agreement**

Sensor readings should exhibit strong linear correlations with their corresponding pollutant ground-truth values if functioning properly.

- **H6: Detection of Sensor Faults**

A significant negative or inconsistent correlation between a sensor and its reference value indicates possible sensor drift or miscalibration.

Hypothesis Test Results (1/2)

Hypothesis	Analytical Focus	Statistical Test	Result
H1	Difference in pollutant means between peak and non-peak hours	Independent t-test	Supported
H2	Comparison of pollutant means between winter and non-winter periods	Independent t-test	Supported
H3	Correlation of CO, NO _x with RH; O ₃ with Temperature	Pearson correlation	Not Supported (Weak)

Hypothesis Test Results (2/2)

Hypothesis	Analytical Focus	Statistical Test	Result
H4	Variation in benzene (C ₆ H ₆) across weekdays	One-way ANOVA	Supported
H5	Correlation between sensor and ground-truth values	Pearson correlation	Partially Supported
H6	Detection of sensor anomaly (NO _x sensor)	Pearson correlation	Supported

Key Findings

This study provided a data-driven investigation into urban air pollution using real-world sensor data.

Key conclusions include:

- **Pollutant levels peak during rush hours and winter months**, increasing health risks for urban populations.
- **CO and Benzene sensors performed reliably**, while the **NOx sensor showed significant calibration issues**.
- Weather conditions such as temperature and humidity influence pollutant levels, though the correlations are modest.
- **Statistical testing confirmed key patterns**, supporting targeted interventions in urban traffic and sensor monitoring.

Overall, this project demonstrates how environmental data can be used to guide public health policies and smarter city planning.

Public Health Implications

Our analysis highlights specific time periods and environmental conditions where air quality poses elevated health risks:

- **Weekday rush hours** (8 AM and 6 PM) show sharp increases in CO and NOx levels
- **Winter months** exhibit sustained high pollution due to heating and stagnant air
- These patterns raise concerns for:
 - **Children**, who may commute during peak hours
 - **Elderly individuals** and **asthma patients**, who are more vulnerable to pollutants
- **Unreliable sensors**, especially for NOx, pose a risk for delayed response and misinformed policy decisions

Policy Recommendations

Based on our findings, we propose the following interventions to improve urban air quality and protect public health:

- **Traffic Management**
 - Implement congestion pricing or restrict traffic during peak hours
 - Promote public transportation and carpooling incentives
- **Urban Planning**
 - Avoid placing schools and hospitals near major roads
 - Enforce zoning laws that consider air quality patterns
- **Sensor Maintenance**
 - Establish calibration protocols for all air quality sensors, particularly for NO_x
 - Deploy multiple sensors per location to detect faults early
- **Public Awareness**
 - Promote awareness campaigns around peak pollution times and seasons

Future Work

To extend the impact of this study, future research can explore:

- **Spatio-temporal analysis** using data from multiple locations or mobile sensors
- **Integration with healthcare data** to study links between pollution exposure and hospital admissions
- **Machine learning models** to predict pollutant surges and detect sensor anomalies
- **City-scale pollution mapping** with GIS tools and real-time feeds
- **Simulation of policy scenarios** to estimate the effect of interventions like car bans or zoning reforms

Thank You!

Group 28 – Spring 2025