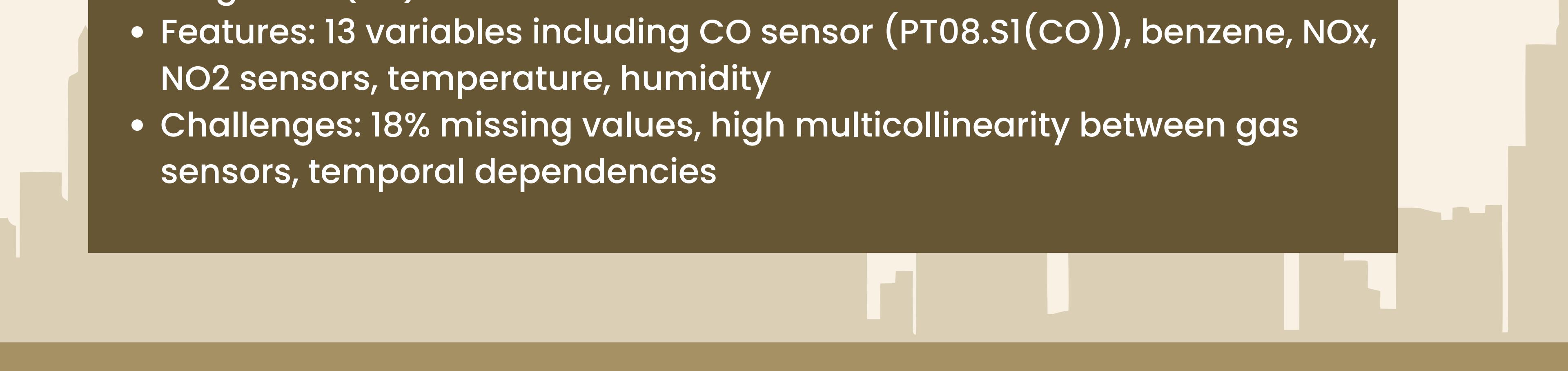


# AIR QUALITY ANALYSIS

## Members :

- **Simran Amesar**
- **Shreya Chowdary Challagulla**
- **Ishwari Thakur**
- **Pooja Chandrappa**
- **Yuvraj Ghag**





# DATASET



The UCI Air Quality Dataset contains 9,357 hourly observations collected over 13 months (Mar 2004 - April 2005) from a single urban monitoring station.

- Target: CO(GT) - true CO concentration
- Features: 13 variables including CO sensor (PT08.S1(co)), benzene, NOx, NO<sub>2</sub> sensors, temperature, humidity
- Challenges: 18% missing values, high multicollinearity between gas sensors, temporal dependencies

# PREPROCESSING



- Missing Values: Values coded as -200 represent sensor failures
- Total missing: 16,701 instances across 13 numeric features (18% of data)
- Treatment: Median imputation (robust to skewness)
- Temporal Structure: Hourly measurements with daily/seasonal patterns
- Date/Time columns combined into DateTime index
- Target-Feature Separation : 'CO(GT)'

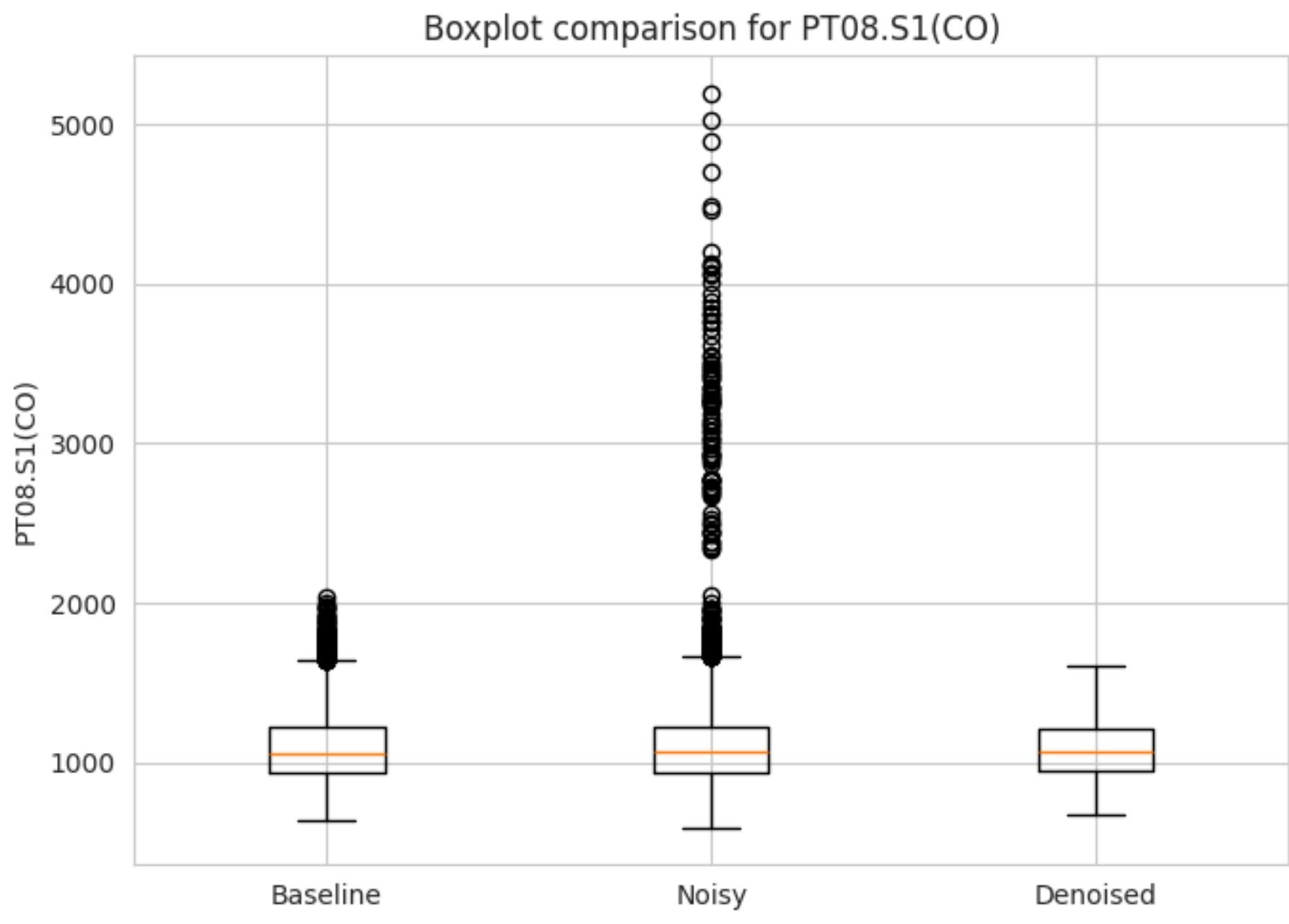
# NOISE INJECTION

- Add artificial noise to 10% of samples to test data cleaning methods
- Randomly inject outlier values (0-1 range) into selected features : "PT08.S1(CO)" and "PT08.S2(NMHC)"
- Tests outlier detection performance and model robustness against noisy data

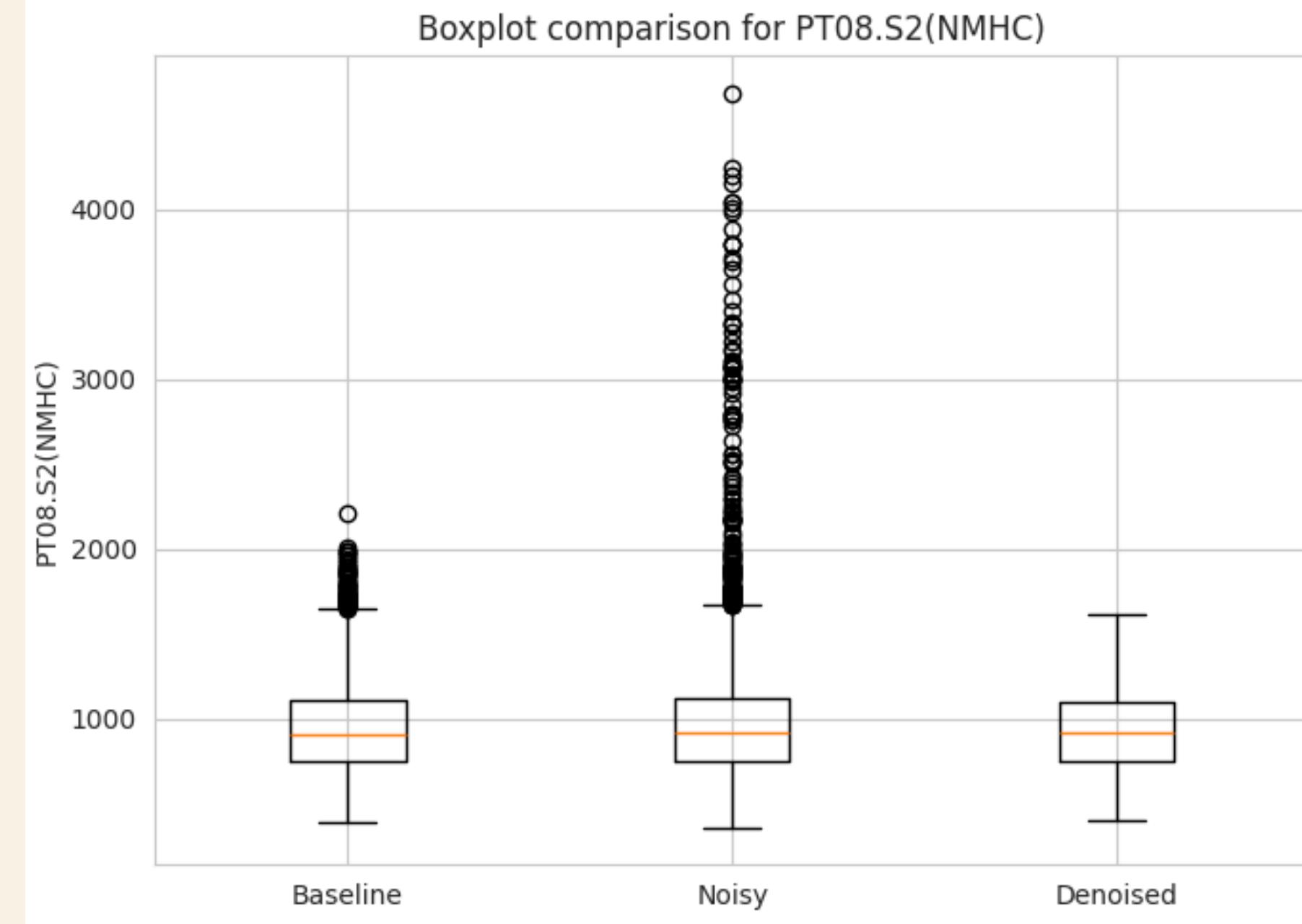


# NOISE INJECTION

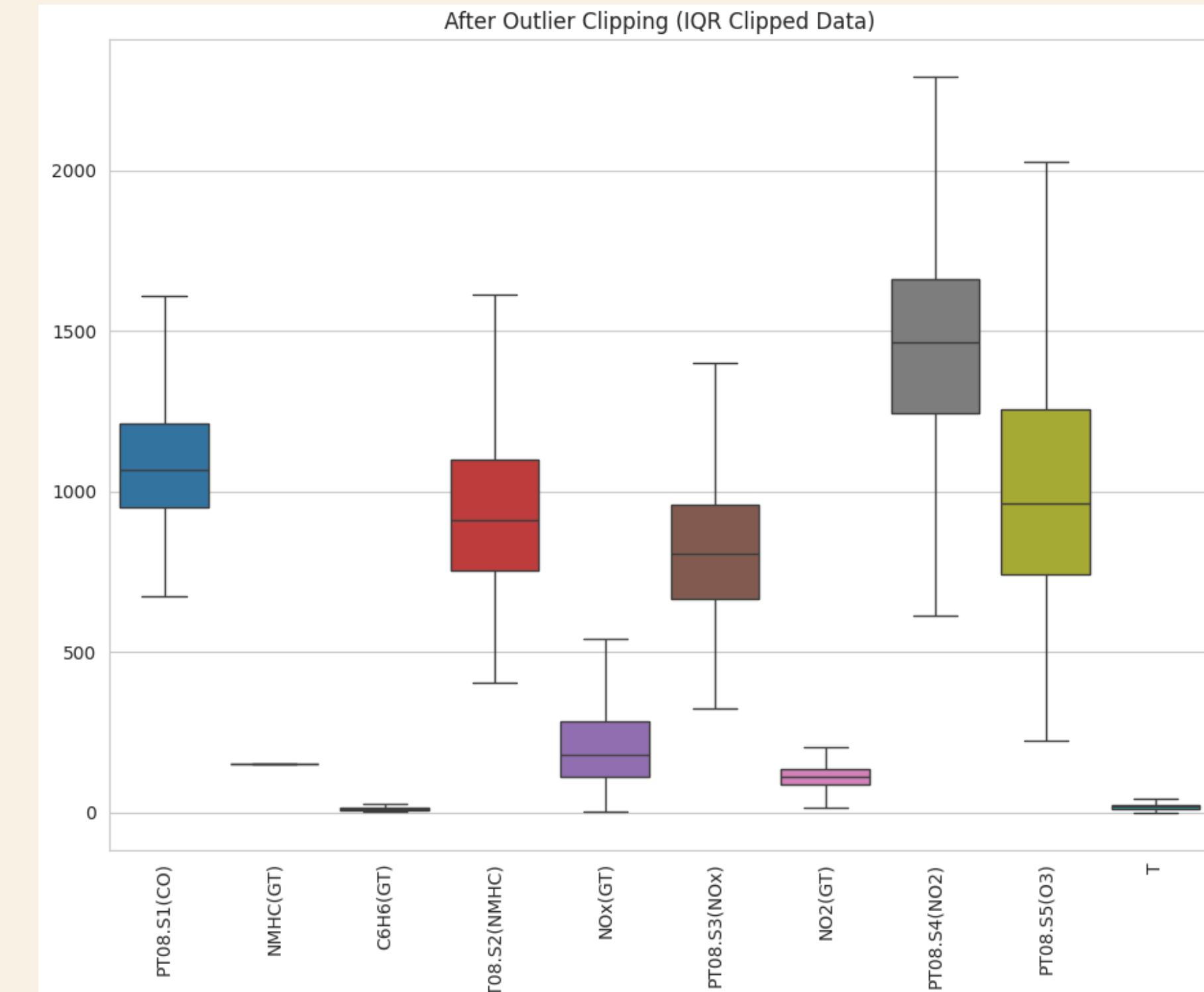
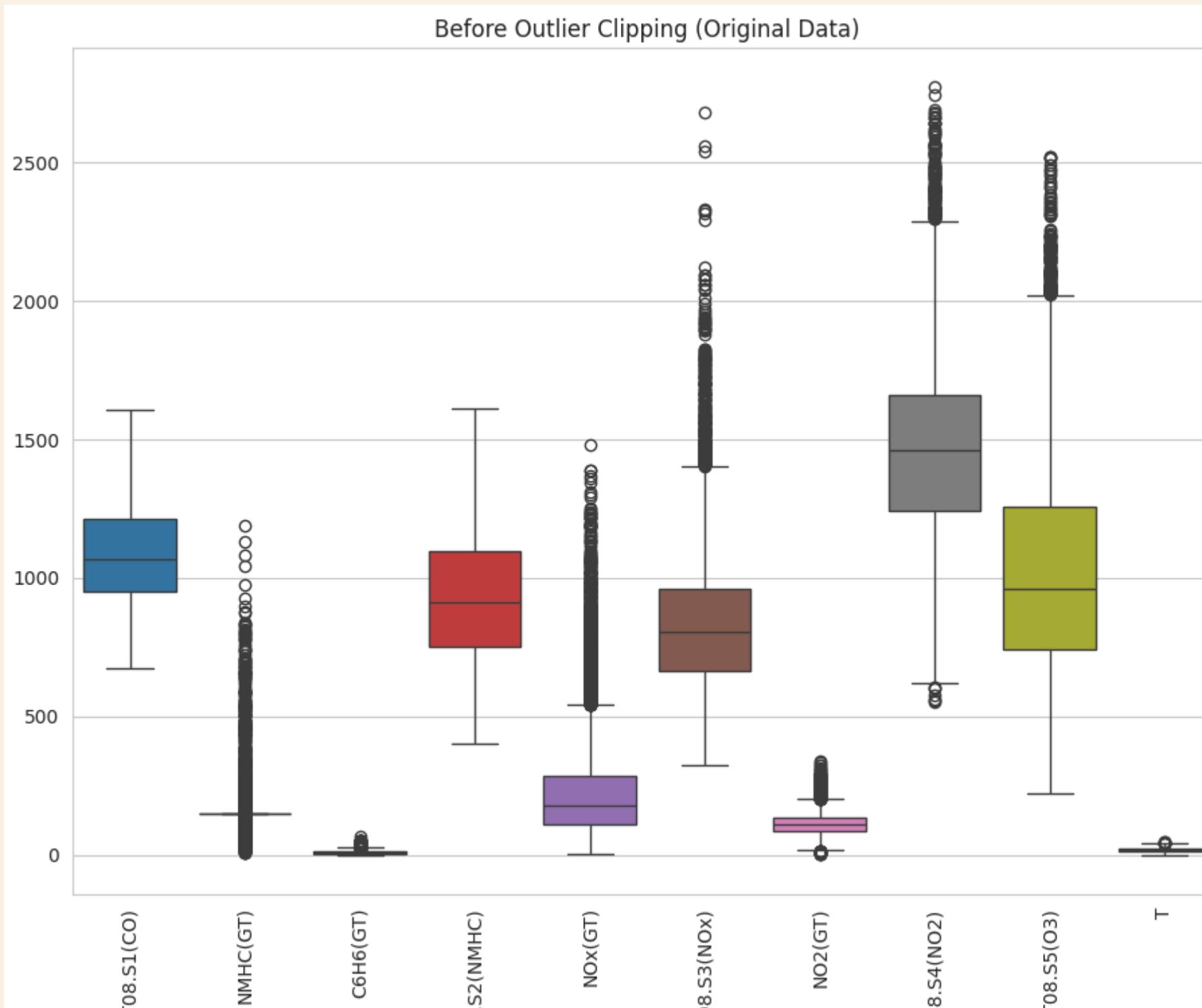
Boxplot comparison for PT08.S1(CO)



Boxplot comparison for PT08.S2(NMHC)



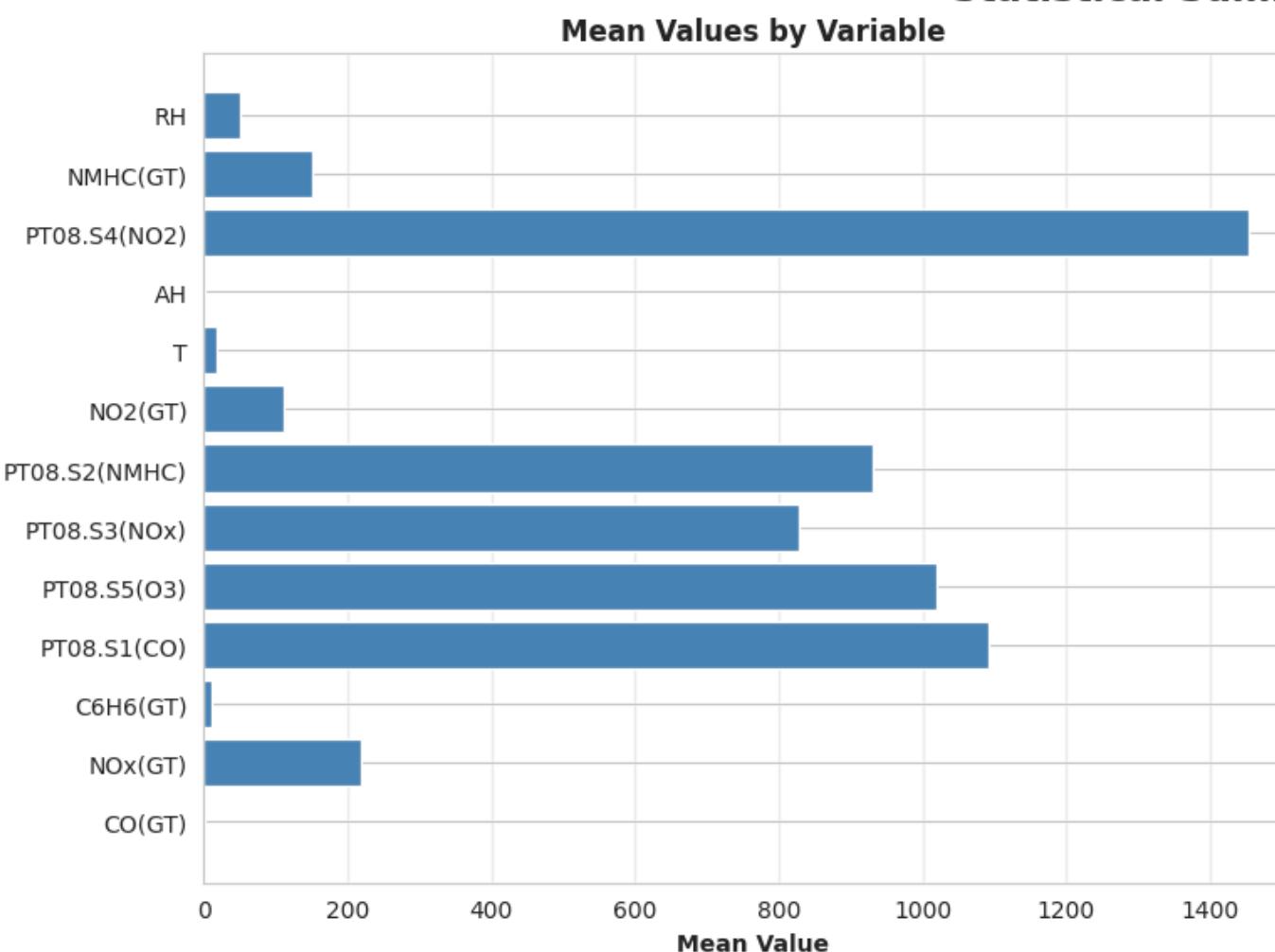
# OUTLIER CLIPPING



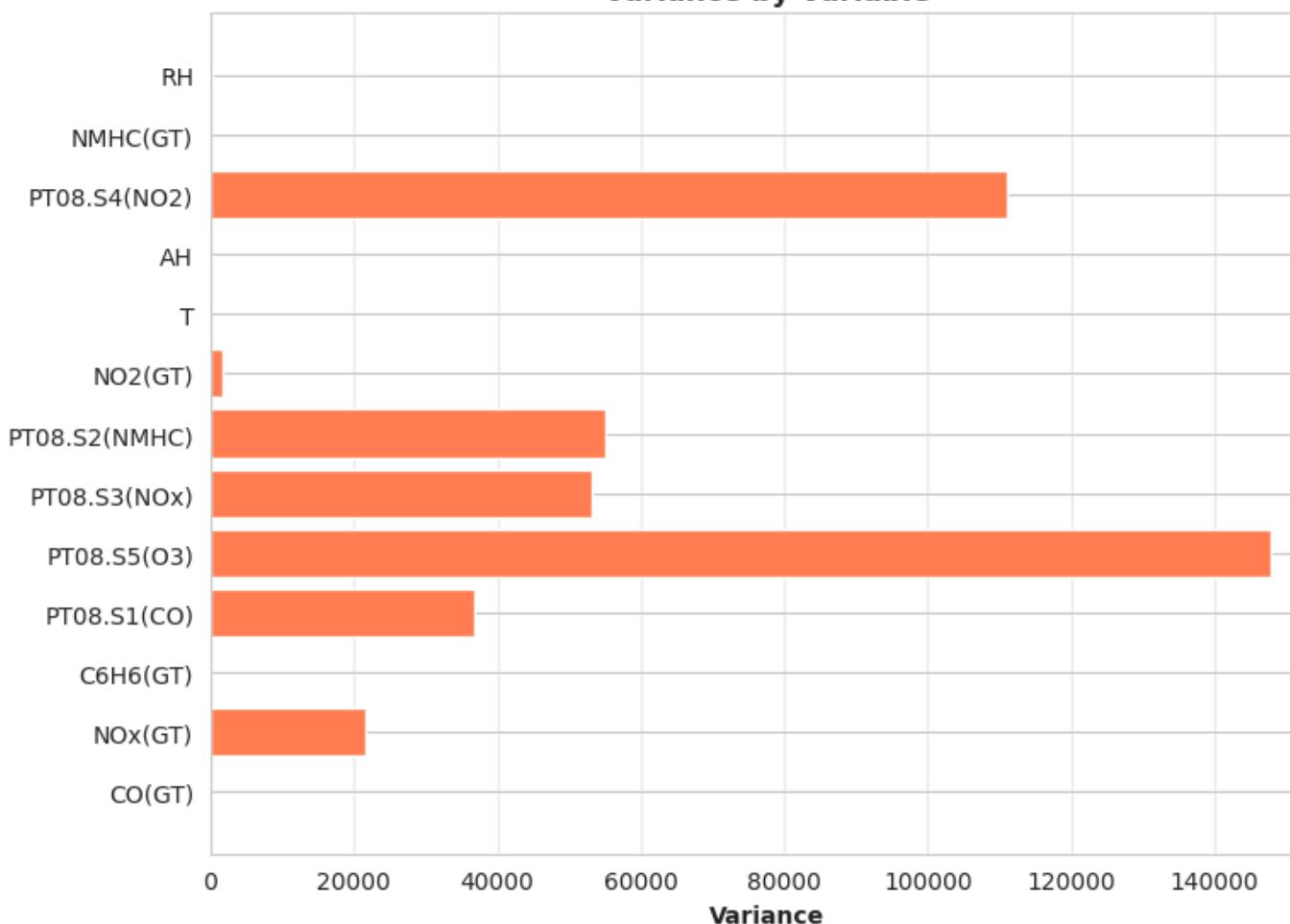
# STATISTICAL ANALYSIS



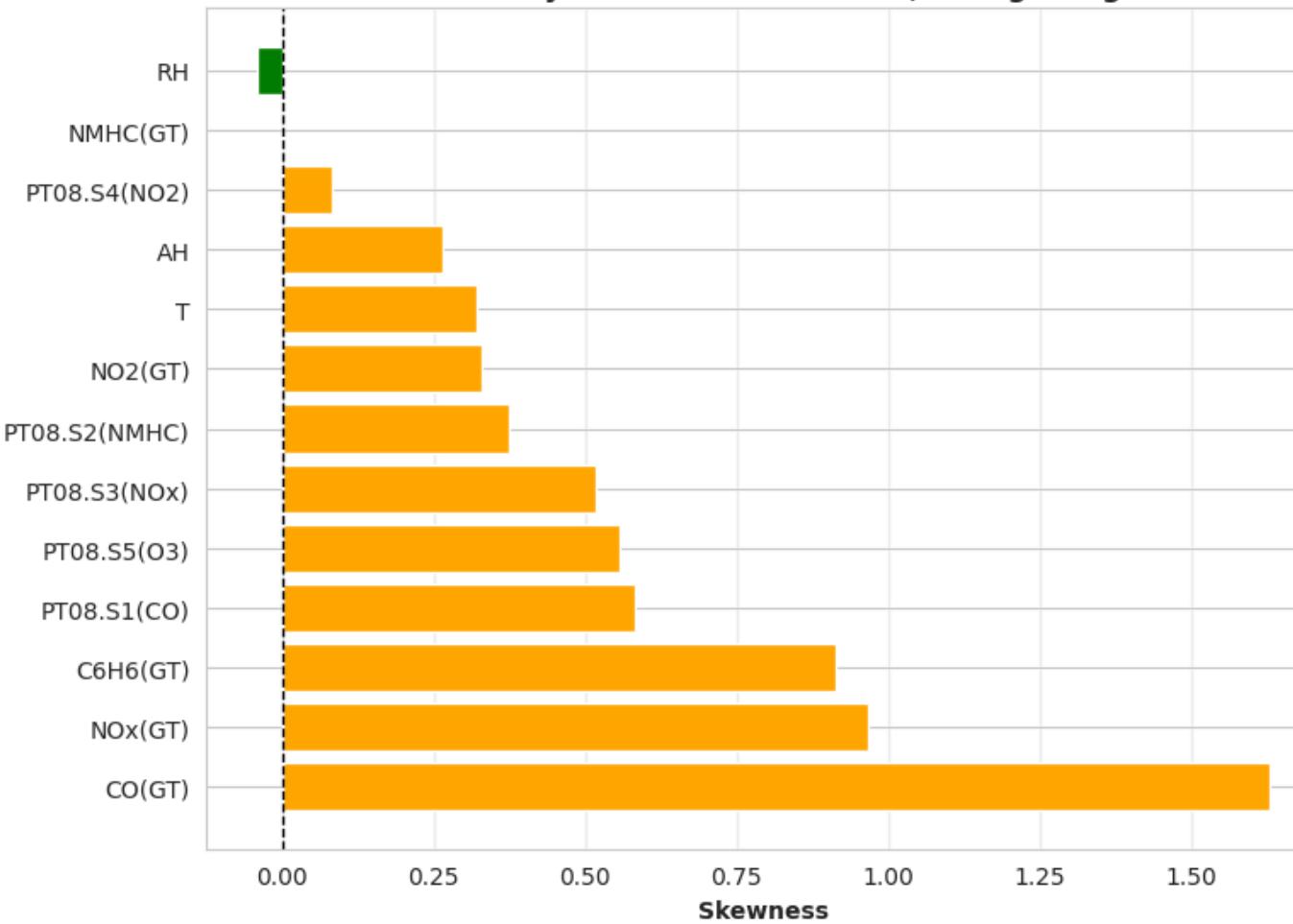
## Statistical Summary - All Measures



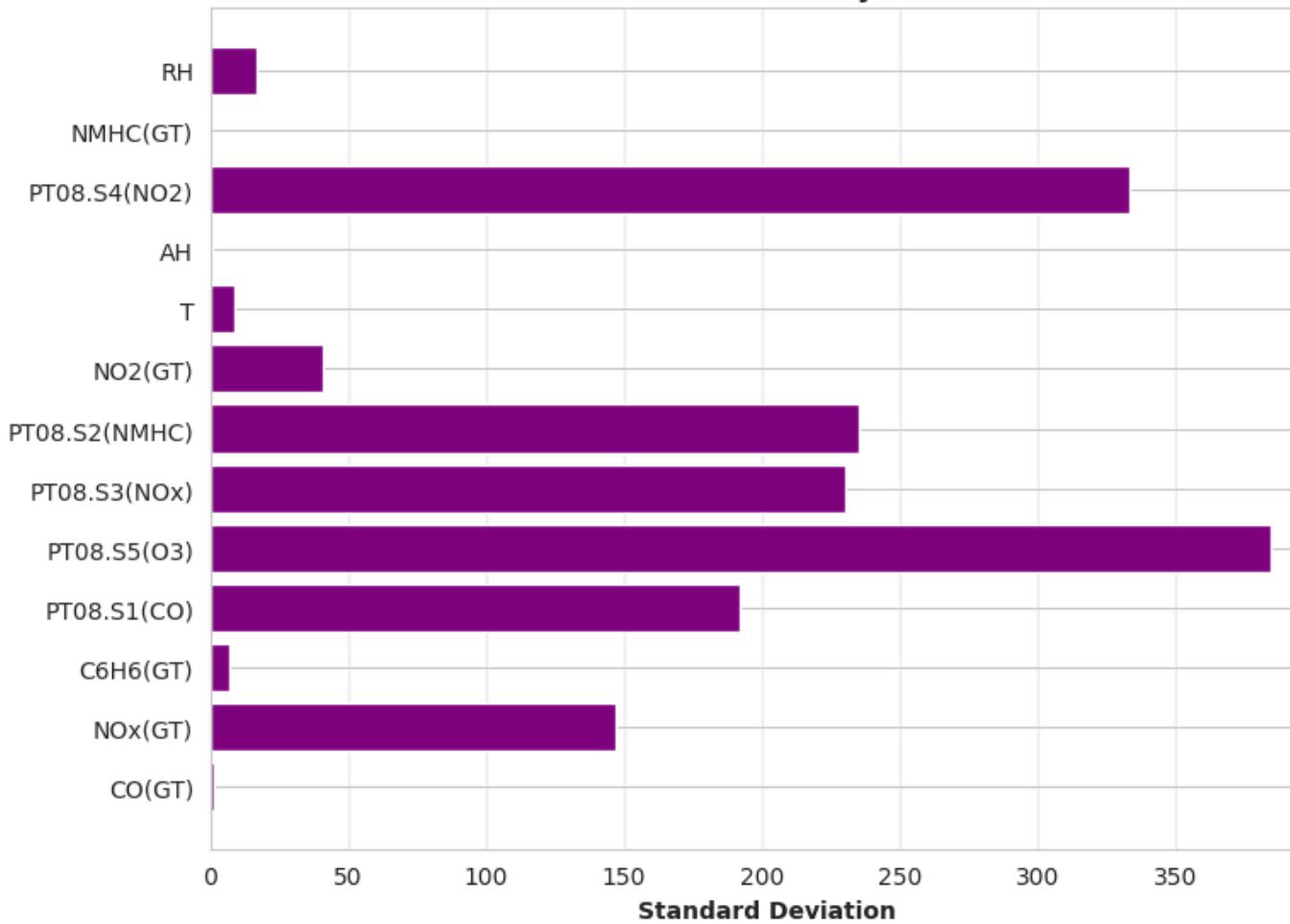
Variance by Variable

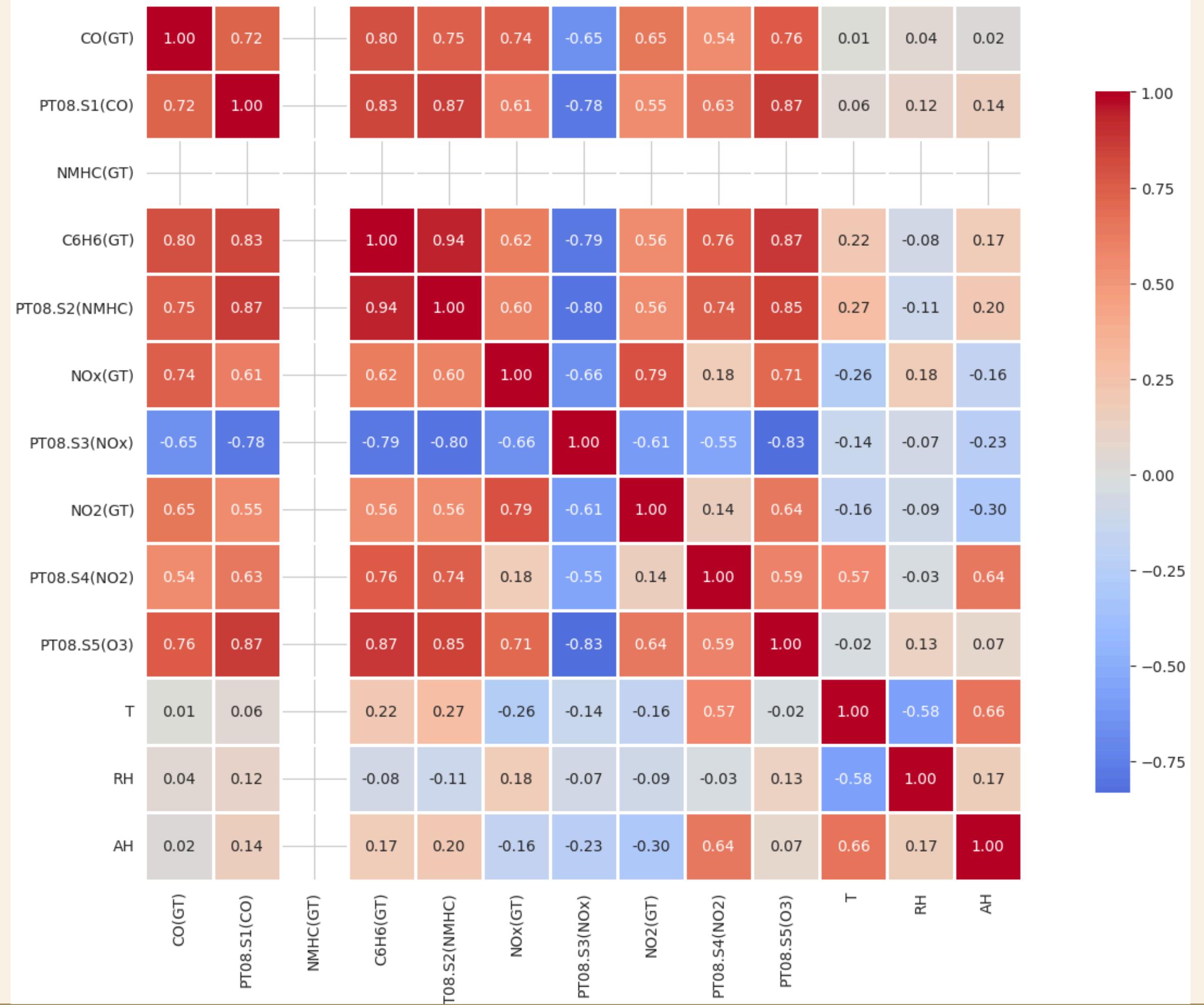


Skewness by Variable (Green=Left, Orange=Right)



Standard Deviation by Variable



**Correlation Heatmap - Air Quality Variables**

# **SELECTED FEATURES**

**PT08.S1(CO)**

**C6H6(GT)**

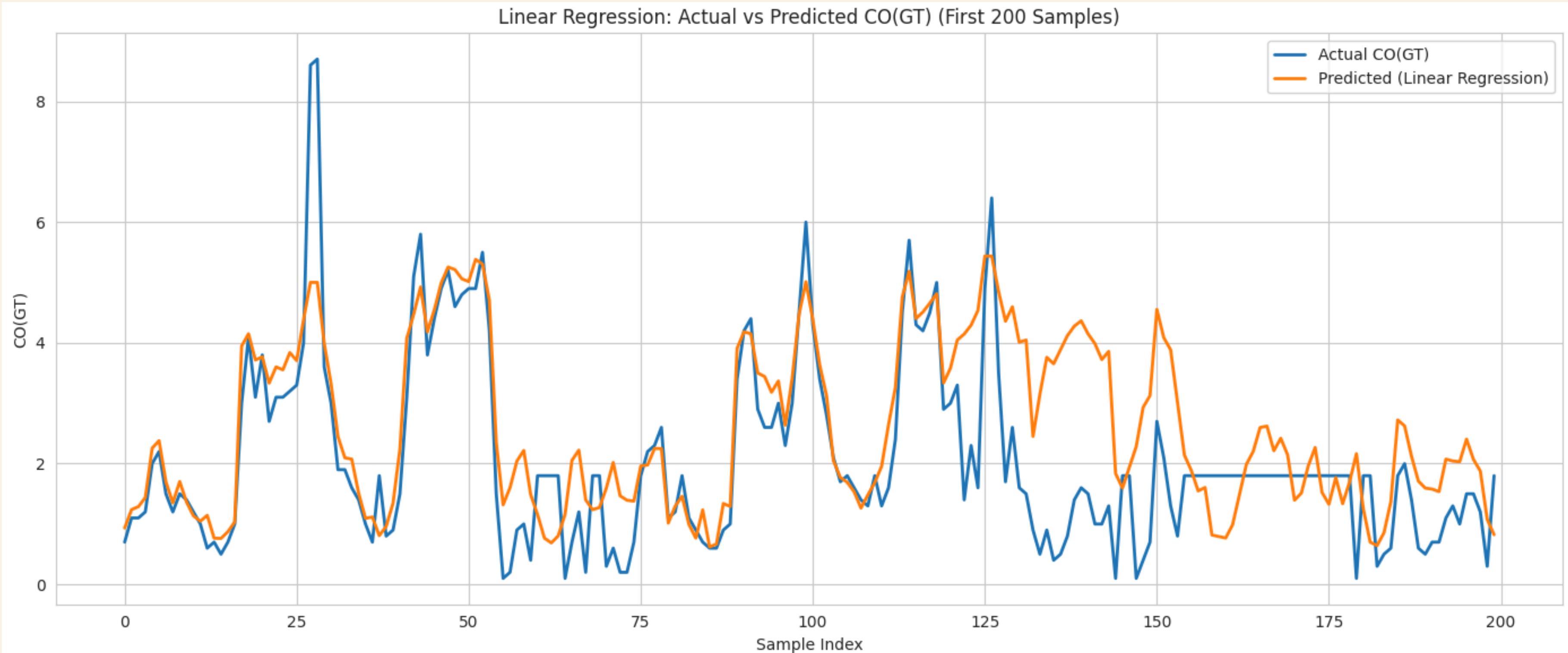
**PT08.S2(NMHC)**

**NOx(GT)**

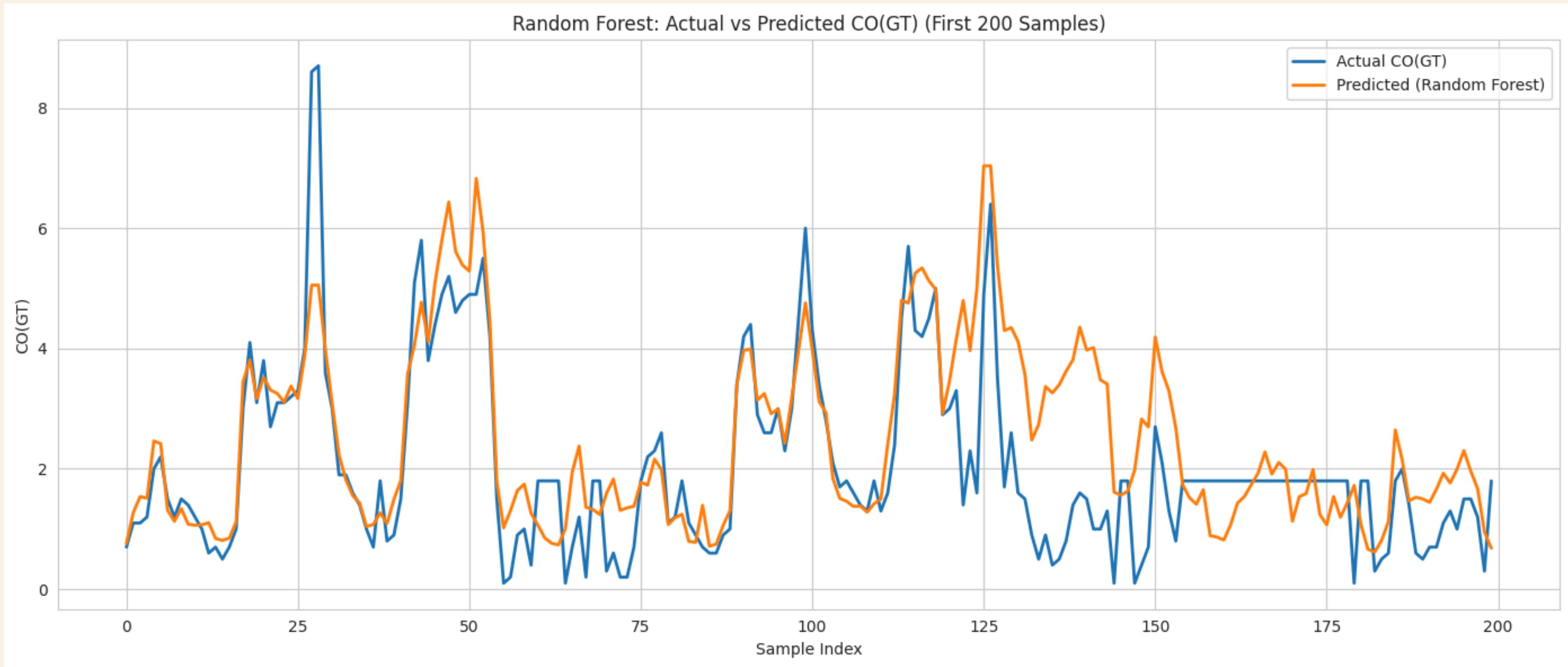
**PT08.S5(O3)**

# **SUPERVISED LEARNING**

# LINEAR REGRESSION



# RANDOM FOREST

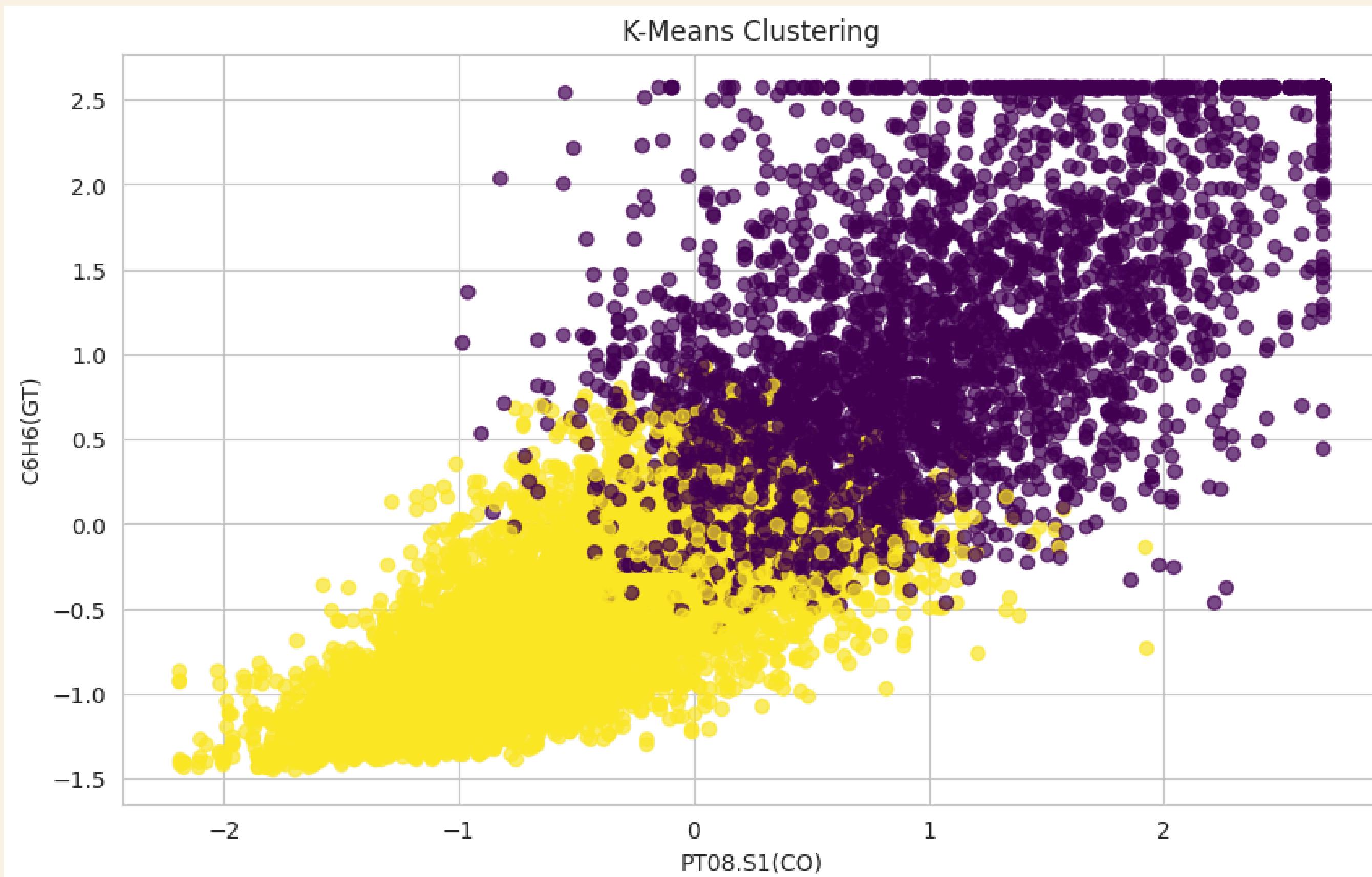


# MODEL COMPARISON

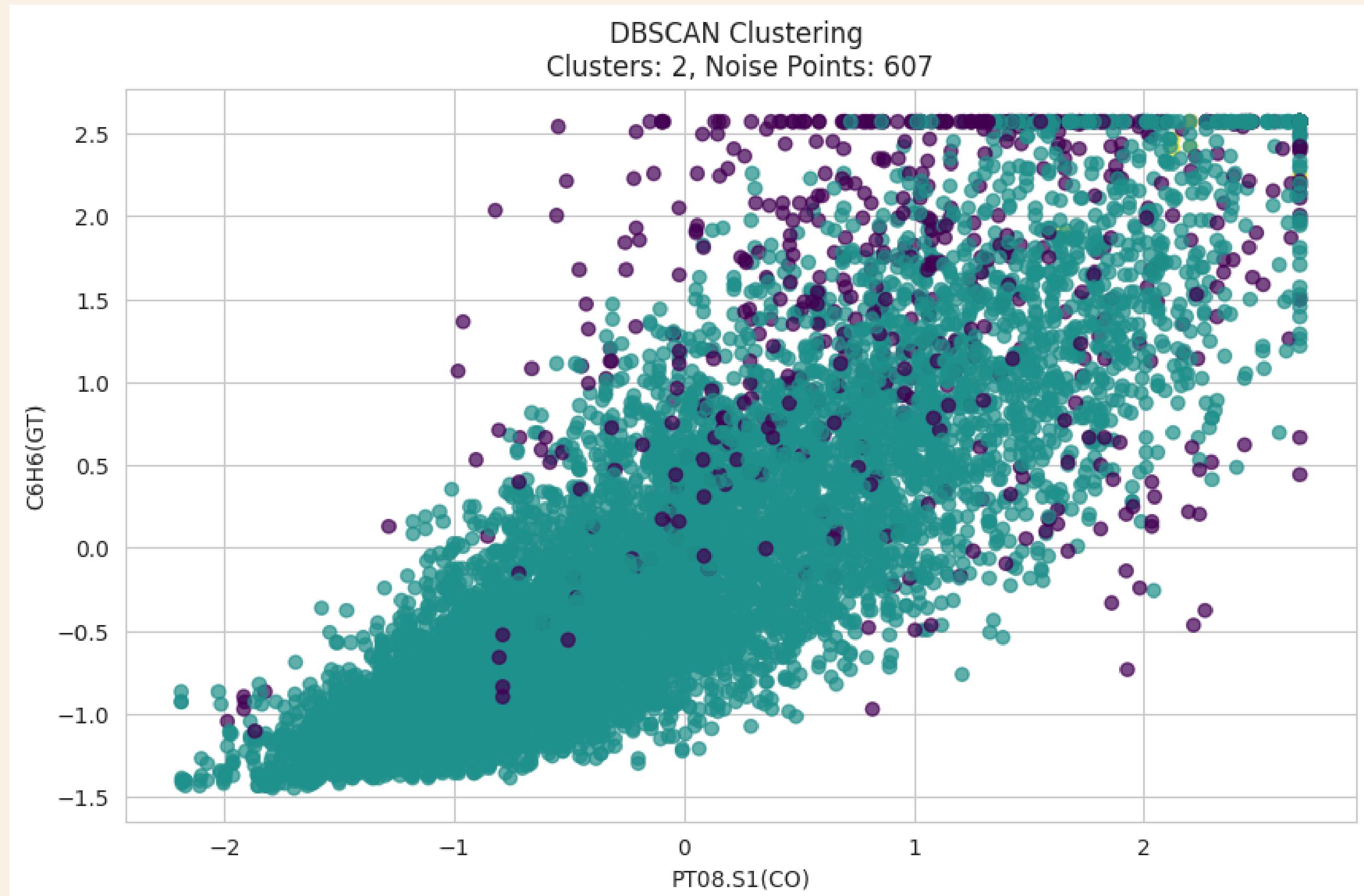
Model	RMSE	MAE	R2
Random Forest	0.655281	0.435162	0.760967
Linear Regression	0.726689	0.520099	0.706032

# **UNSUPERVISED LEARNING**

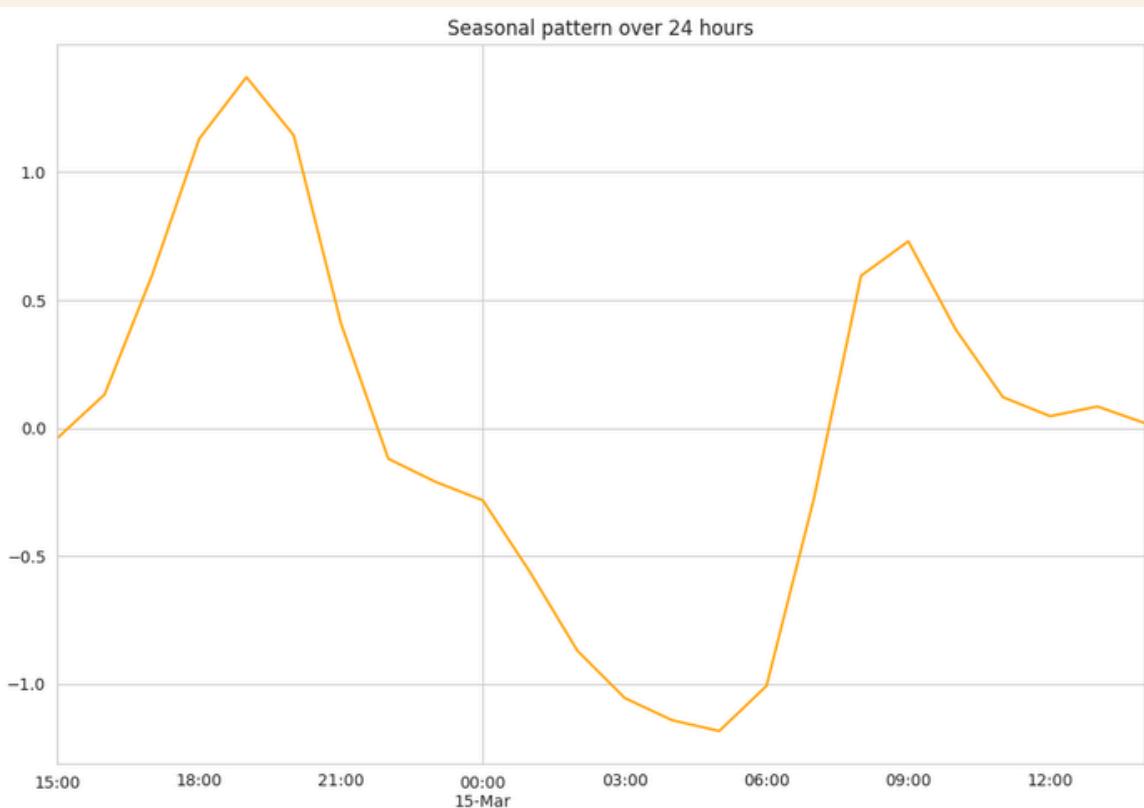
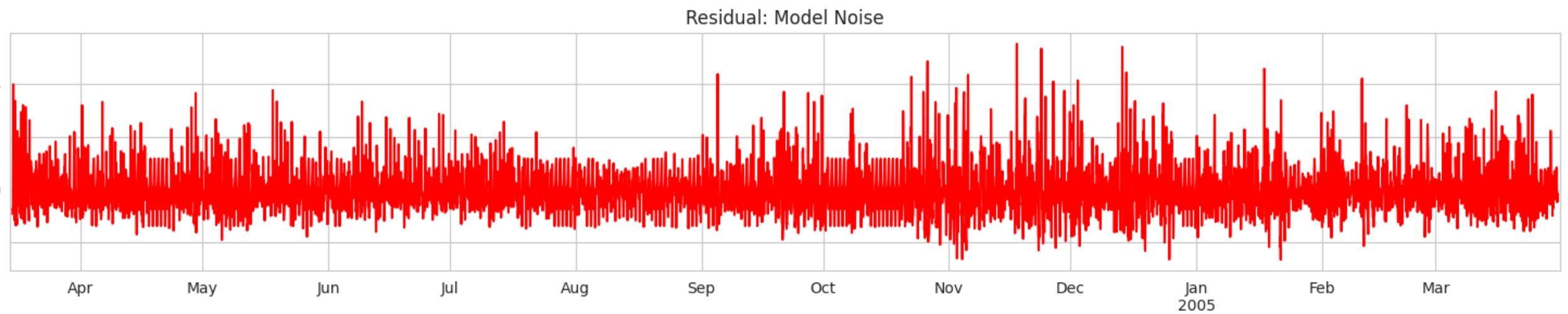
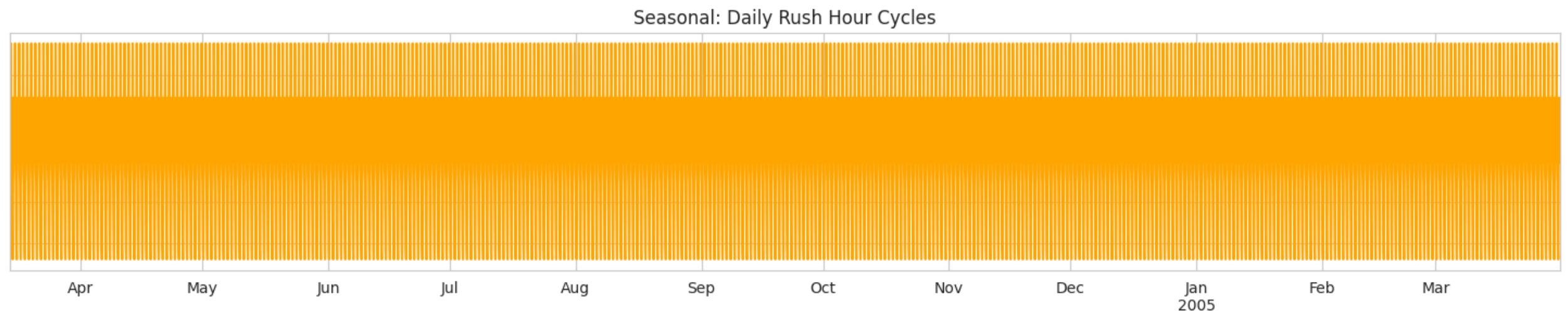
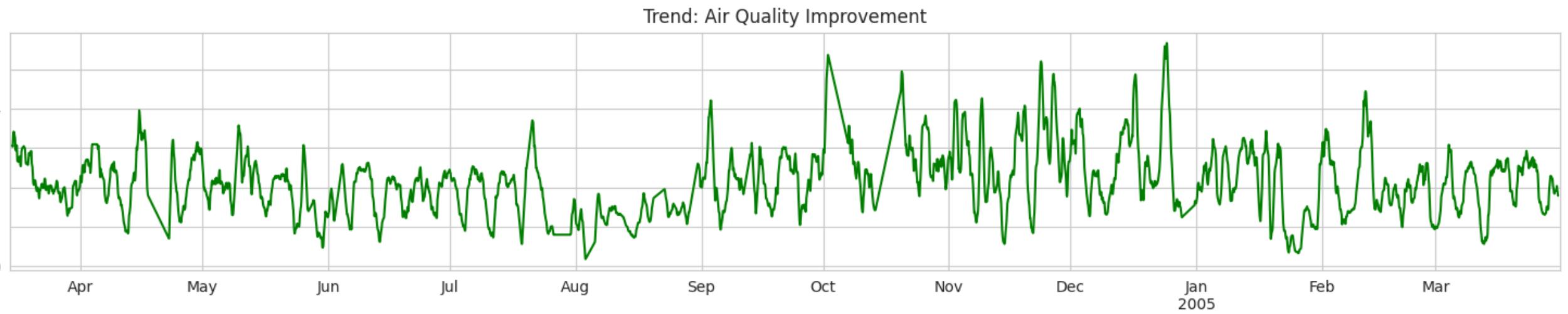
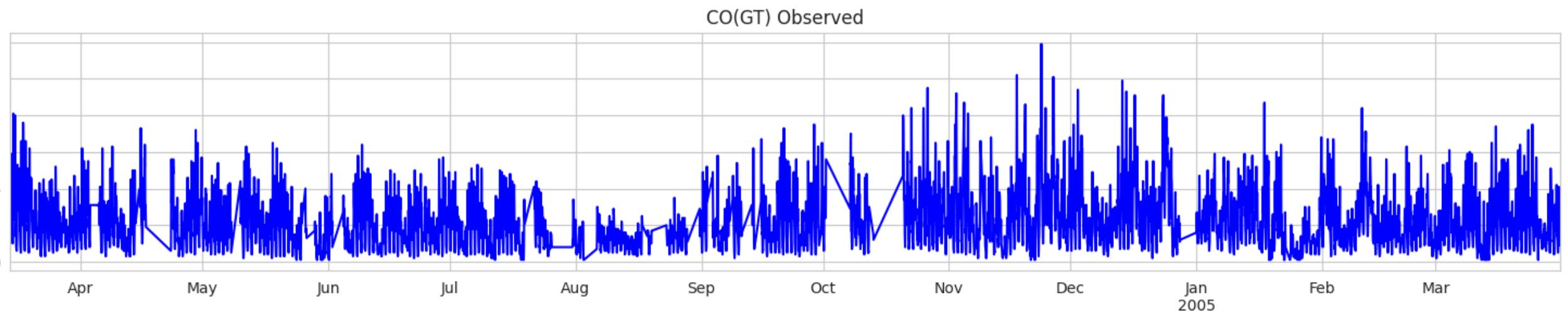
# K-MEANS CLUSTERING



# DBSCAN CLUSTERING



# TIME-SERIES DECOMPOSITION



# THANK YOU

