# Air Quality

By: Krish Shah

# Hypothesis

**My hypothesis is that there will be a positive correlation between higher temperatures and lower humidity with elevated concentrations of pollutants. In this study pollutants include CO, NOx, and C6H6. This correlation indicates a link between weather conditions and the degradation of air quality.**

Null Hypothesis: There is no significant relationship between environmental factors (temperature, relative humidity, absolute humidity) and the concentrations of air pollutants (e.g., CO, NOx, C6H6).

Alternative Hypothesis: There is a significant relationship between environmental factors (temperature, relative humidity, absolute humidity) and the concentrations of air pollutants (e.g., CO, NOx, C6H6).
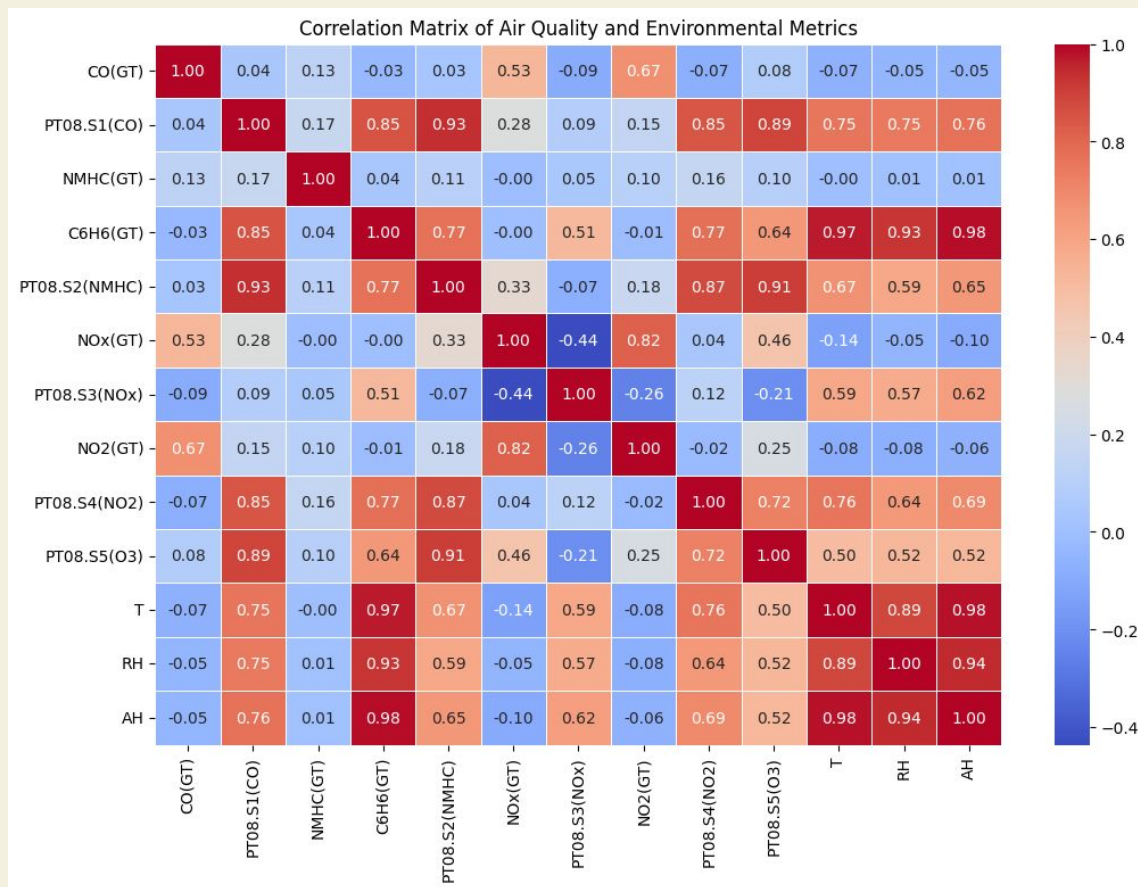
# Column Names

Not all column descriptions are listed.

- Date – Date of data recording.
- Time – Time of data recording.
- CO(GT) – Carbon monoxide concentration in mg/m³ (ground truth).
- PT08.S1(CO) – Sensor reading for CO concentration (indirect measurement).
- NMHC(GT) – Non-methane hydrocarbons concentration in µg/m³ (ground truth).
- C6H6(GT) – Benzene concentration in µg/m³ (ground truth).
- PT08.S2(NMHC) – Sensor reading for NMHC concentration (indirect measurement).
- NOx(GT) – Nitrogen oxides concentration in ppb (ground truth).
- PT08.S3(NOx) – Sensor reading for NOx concentration (indirect measurement).
- NO2(GT) – Nitrogen dioxide concentration in µg/m³ (ground truth).
- PT08.S4(NO2) – Sensor reading for NO2 concentration (indirect measurement).
- PT08.S5(O3) – Sensor reading for ozone (O3).
- T – Ambient temperature in Celsius.
- RH – Relative humidity in percentage.
- AH – Absolute humidity in g/m³.

# Heat Map

There is clear visualization of the correlation between all the variables from each column



Correlation Matrix of Air Quality and Environmental Metrics

# Linear Regression

The data has been split into training sets to be used for linear regression. In this list only the intercept is listed. The coefficients show the correlation between the independent variables(T, RH, AH) and dependent variables(pollutant concentrations).

**Coefficients(Intercept)**:

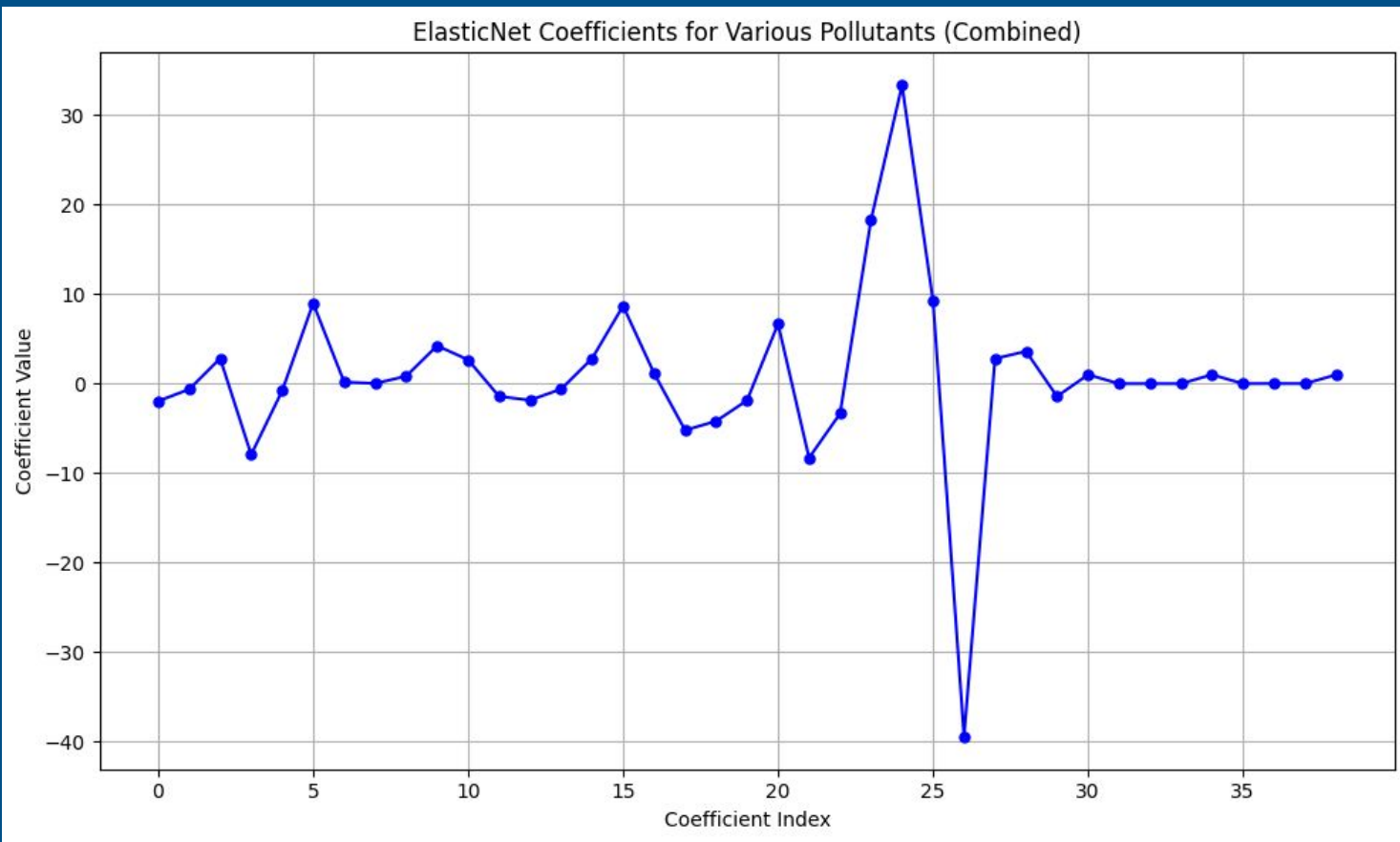| | |
|---|---|
| CO(GT) | 31.0791 |
| PT08.S1(CO) | 890.7222 |
| NMHC(GT) | -96.5638 |
| C6H6(GT) | 5.4098 |
| PT08.S2(NMHC) | 724.3954 |
| NOx(GT) | 343.7566 |
| PT08.S3(NOx) | 1142.6391 |
| NO2(GT) | 221.9128 |
| PT08.S4(NO2) | 409.1585 |
| PT08.S5(O3) | 793.7094 |
| T | $5.3291e{-15}$ |
| RH | $-4.9738e{-14}$ |
| AH | $-1.3323e{-14}$ |

# Elastic Net Coefficients

These intercepts represent the baseline values of the pollutants when all the independent variables (T, RH, AH) are set to zero in the Elastic Net regression model. **Lasso** identifies the most important predictors by simplifying the model. **Ridge** reduces the size of coefficients without setting any to zero, allowing for relevant values to stay in the model.

Lasso and Ridge were used in combination, which allowed for the exclusion of unnecessary features and prevention of overfitting. My model shows these results in a manner that balances both feature selection and coefficient shrinkage.
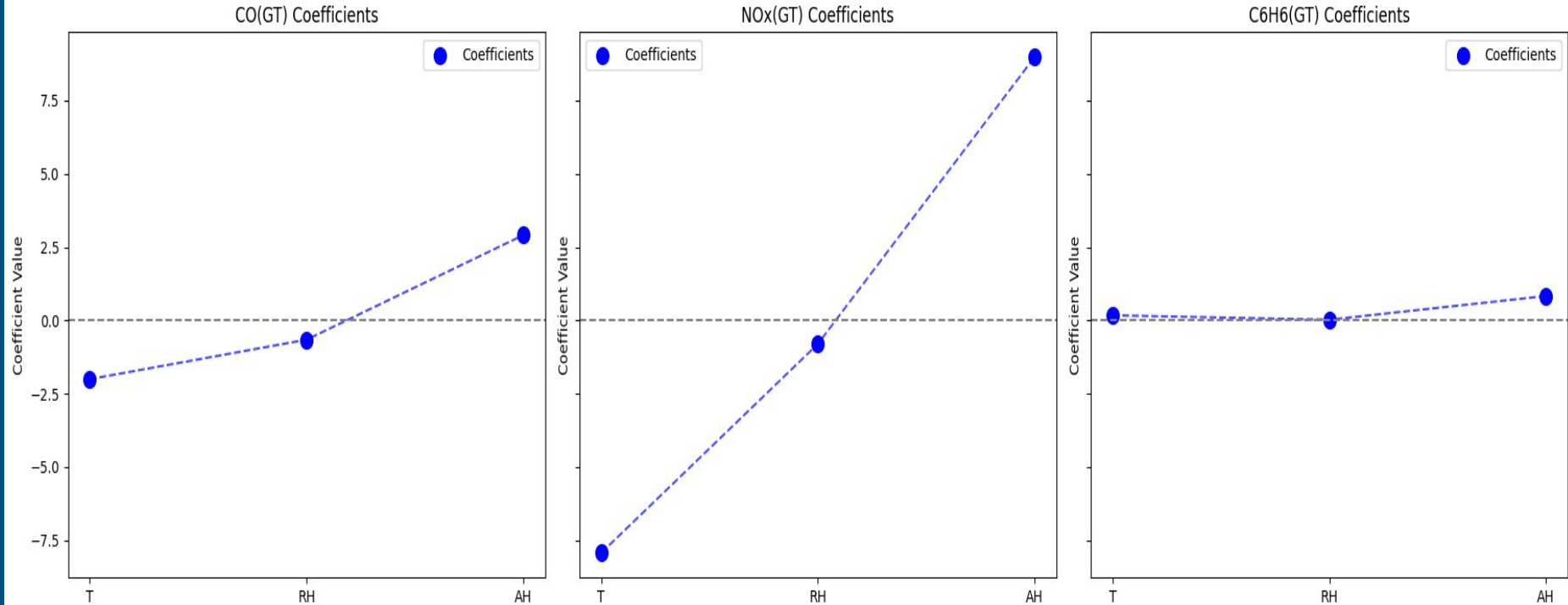
CO(GT): 29.093
NOx(GT): 338.590
C6H6(GT): 4.997
PT08.S1(CO): 892.463
NMHC(GT): −98.525
PT08.S2(NMHC): 728.490
NO2(GT): 218.104
PT08.S3(NOx): 1134.066
PT08.S4(NO2): 430.031
PT08.S5(O3): 795.272
T: −0.005
RH: 0.009
AH: −0.516

# Elastic Net Coefficients Graph



ElasticNet Coefficients for Various Pollutants (Combined)

# Dot Plot Subplots



Elastic Net Coefficients Compared to Independent Variables

# Thank You

Air Quality Correlation