

Milestone 4

Anh Le, Eve Wicksteed

The impact of weather on air quality

I. Introduction

The adverse affects of air pollution on health are well documented and air pollution can lead to a large range of diseases and increased morbidity and mortality (Younger et al., 2008). Adverse health impacts include, but are not limited to, lung cancer risk, respiratory infections, allergic disease and asthma (Younger et al., 2008; Shea et al., 2008). These health risks can affect a large proportion of the population as many different groups are vulnerable to the effects of air pollution including infants, children, the elderly, people with impaired immune systems, and people who work or are physically active outdoors (Matooane et al., 2004).

Because of the many, and severe, impacts of air quality, it is important to understand patterns in the data. We have a dataset of air quality observations as well as temperature and humidity data which we will use to gain understanding of the patterns and impacts of weather on air quality.

II. Research question

As stated above we would like to understand the impact of weather on air quality. For this reason our research question is: - What is the affect of temperature and humidity on the concentration of air pollutants, such as benzene, titania, and tin oxide?

III. Data Description

The air quality dataset used in this analysis was obtained from the University of California Irvine Machine learning Repository. It was contributed by Saverio De Vito from the National Agency for New Technologies, Energy and Sustainable Economic Development.

The dataset contains 15 variables and 9358 observations of hourly averaged responses from an Air Quality Chemical Multisensor Device. Data were recorded from March 2004 to February 2005, in a significantly polluted area, at road level, within a city in Italy. Variables include the date and time each response was recorded, and the corresponding concentrations of 13 air pollutants analyzed by the sensor device. Missing values are tagged with -200 value. Below is the entire variable set:

Variables	Type	Description
Date	character	Date (DD/MM/YYYY)
Time	time	Time (HH.MM.SS)
CO(GT)	double	True hourly averaged concentration CO in mg/m^3 (reference analyzer)
PT08.S1(CO)	integer	PT08.S1 (tin oxide) hourly averaged sensor response (nominally CO targeted)
NMHC(GT)	integer	True hourly averaged overall Non Metanic HydroCarbons concentration in microg/m^3 (reference analyzer)
C6H6(GT)	double	True hourly averaged Benzene concentration in microg/m^3 (reference analyzer)
PT08.S2(NMHC)	integer	PT08.S2 (titania) hourly averaged sensor response (nominally NMHC targeted)
NOx(GT)	integer	True hourly averaged NOx concentration in ppb (reference analyzer)
PT08.S3(NOx)	integer	PT08.S3 (tungsten oxide) hourly averaged sensor response (nominally NOx targeted)
NO2(GT)	integer	True hourly averaged NO2 concentration in microg/m^3 (reference analyzer)
PT08.S4(NO2)	integer	PT08.S4 (tungsten oxide) hourly averaged sensor response (nominally NO2 targeted)
PT08.S5(O3)	integer	PT08.S5 (indium oxide) hourly averaged sensor response (nominally O3 targeted)

Variables	Type	Description
T	double	Temperature in °C
RH	double	Relative Humidity (%)
AH	double	AH Absolute Humidity

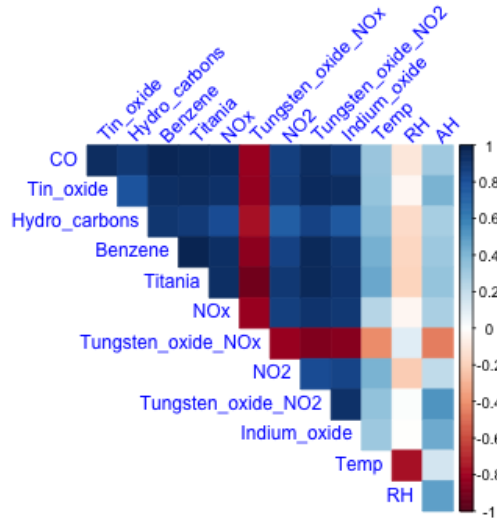
IV. Methods

We are interested in the hourly averaged concentrations of air pollutants, temperature and humidity. We ignore variables which have too many missing data to increase the precision of this analysis. The air pollutants that we will focus on are benzene, titania and tin oxide. After dealing with the missing data, we will perform a linear regression analysis using the OLS (ordinary least squares) method. The dependent variables are the daily concentrations of each air pollutants in $\mu\text{g}/\text{m}^3$. The independent variables are temperature and absolute humidity. Relative humidity is not used due to the high correlation to temperature. Coefficients of relevant variables will be plotted with confidence intervals.

V. Results

1. Exploratory data analysis

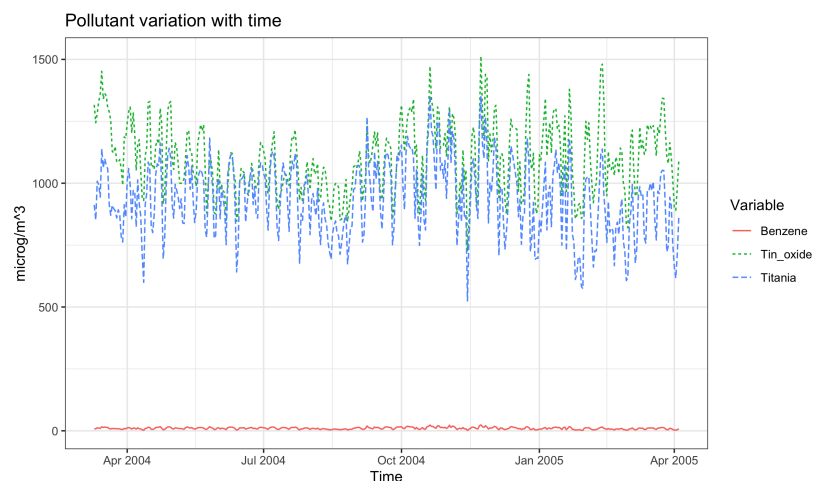
Graph 1: Correlogram of pollutants



Looking at the correlations of the pollutants with weather, we can see that for all pollutants except NOx, temperature (T) is positively correlated, although weakly so. This means that higher temperatures correspond to higher concentrations of the gases. Relative humidity (RH) is negatively and correlated to temperature and has a weak negative correlation to the concentrations of pollutants, except NOx. Absolute humidity (AH) has stronger correlations, mostly positive, although, like temperature, it has a negative correlation with NOx.

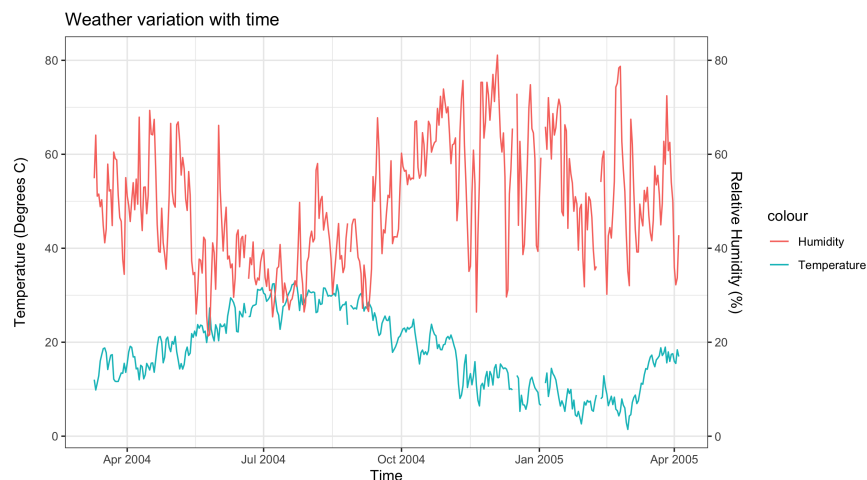
Graph 2: Concentration of some Air Pollutants, Temperature, Humidity over Time, daily average

The plot below shows the **daily** averaged concentrations of some of the pollutants (tin oxide, benzene, and Titania) for a year.



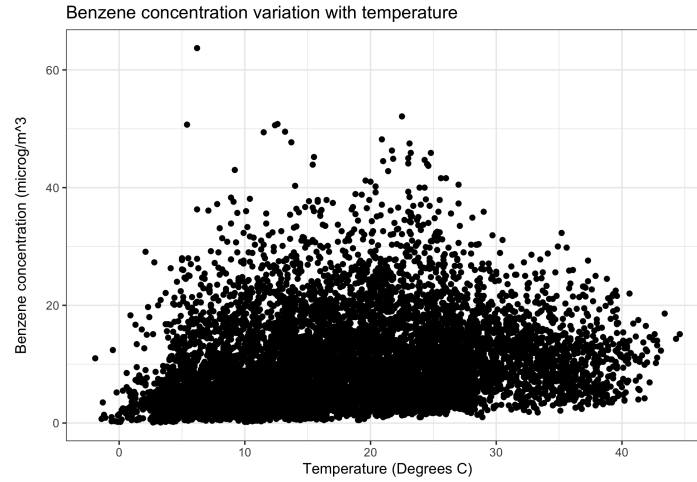
Graph 3: Concentration Temperature and Humidity over Time

The plot below show the **daily** averaged values of temperature and humidity for a year.



Graph 4: Temperature vs. Benzene concentration

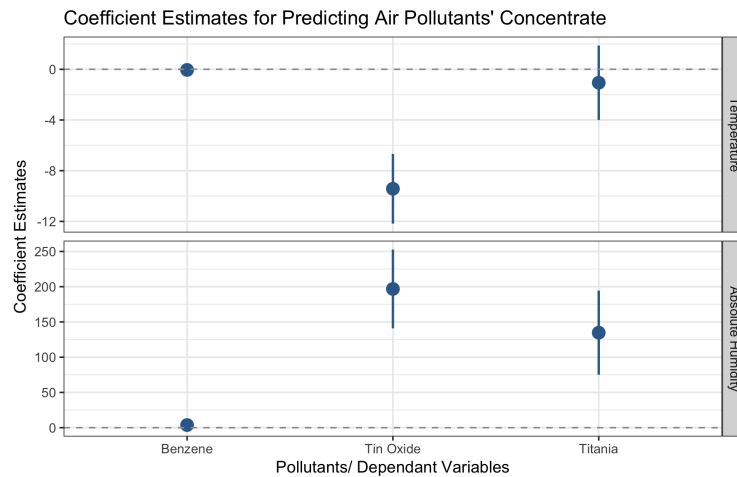
The following graph shows the relationship of benzene to temperature over the year in which data was recorded. The plot suggests there is perhaps a slight relationship. Linear regression in future work will help to clarify the relationships between weather and pollutant concentrations.



2. Linear regression

2.1. Coefficients Plot

We then perform linear regression of all the separate pollutants with temperature and absolute humidity. The graph below shows the coefficients for linear regression for temperature and humidity with all the various pollutants.



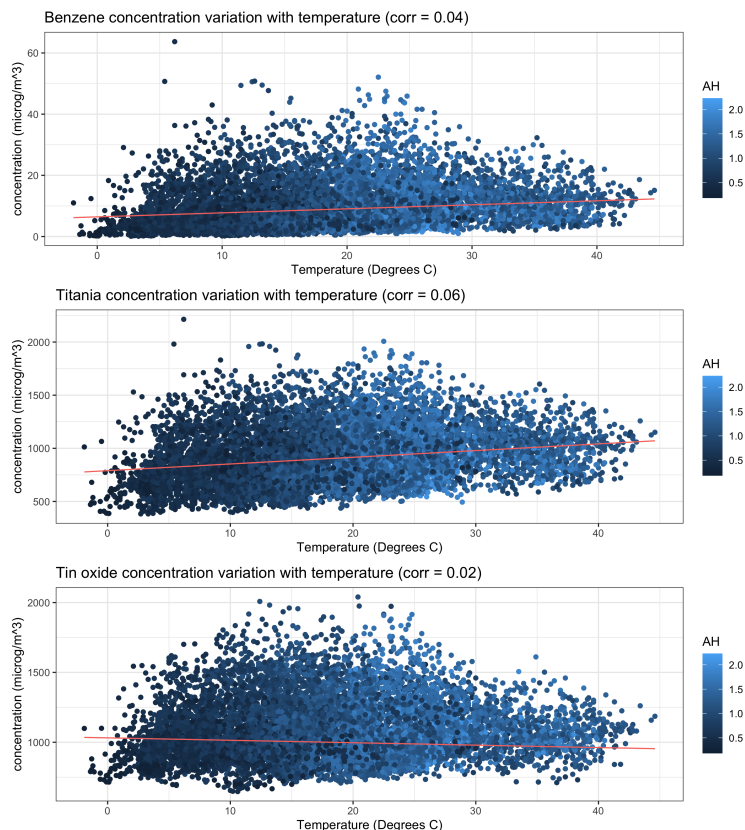
2.2. Regression Output

The table below shows an example of the linear regression output. The dependent variable is the air concentration of Tin oxide in microg/m^3 , while the independent variables are temperature (Temp), and absolute humidity (AH).

```
## # A tibble: 3 x 5
##   term      estimate std.error statistic  p.value
##   <chr>      <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept) 1072.      19.3     55.5 2.82e-184
## 2 Temp       -9.43      1.40    -6.72 6.60e- 11
## 3 AH         197.     28.5     6.90 2.14e- 11
```

VI. Discussion and Conclusion

The next plot shows three of the pollutants (Benzene, Titania and Tin Oxide) plotted with temperature with the linear regression line also plotted. From looking at the plots we can tell the the linear regression line does not represent a lot of the data well. This is shown in the low correlations values (in the plot titles).



It is clear that just looking at the values of temperature and humidity on their own do not provide sufficient information to explain or predict the given concentrations of pollutants. This can be seen in the figures produced above where the linear regression line does not capture well the variation of pollutants. Although weather will affect some pollutants, a more important determinant of pollutant concentrations is emissions. These often vary with time of day and human activity and thus a model incorporating weather as well as these other important variables would likely be more accurate.

References

- S. De Vito, E. Massera, M. Piga, L. Martinotto, G. Di Francia, On field calibration of an electronic nose for benzene estimation in an urban pollution monitoring scenario, *Sensors and Actuators B: Chemical*, Volume 129, Issue 2, 22 February 2008, Pages 750-757, ISSN 0925-4005.
- Matooane, M., John, J., Oosthuizen, R., and Binedell, M. 2004. Vulnerability of South African communities to air pollution. In: 8th World Congress on Environmental Health. Durban, South Africa: Document Transformation Technologies.
- Shea, K., Truckner, R., Weber, R., and Peden, D. 2008. Climate change and allergic disease. *Journal of Allergy and Clinical Immunology*, 122(3): 443-453.
- Younger, M., Morrow-Almeida, H., Vindigni, S., and Dannenberg, A. 2008. The Built Environment, Climate Change, and Health Opportunities for Co-Benefits. *American Journal of Preventative Medicine*, 35 (5):

517-526.