

Predicting Air Quality in Milan

blurb

Dan Herweg

May 2019





Contents

I. Summary

II. Background

I. Motivation

III. Data

I. Methodology

II. Feature Construction

III. Data Exploration

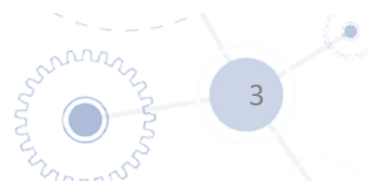
IV. Model

I. Feature Selection

II. Air Quality Index Prediction

III. Bonus: Is smoothed data better?

V. Conclusion



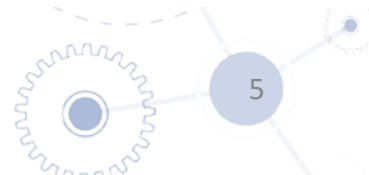
Summary sentence

Summary

- A summary of results
- Pursuing hypothesis that it will work
- This deck is a brief documentation of the project

II. Background

I. Motivation

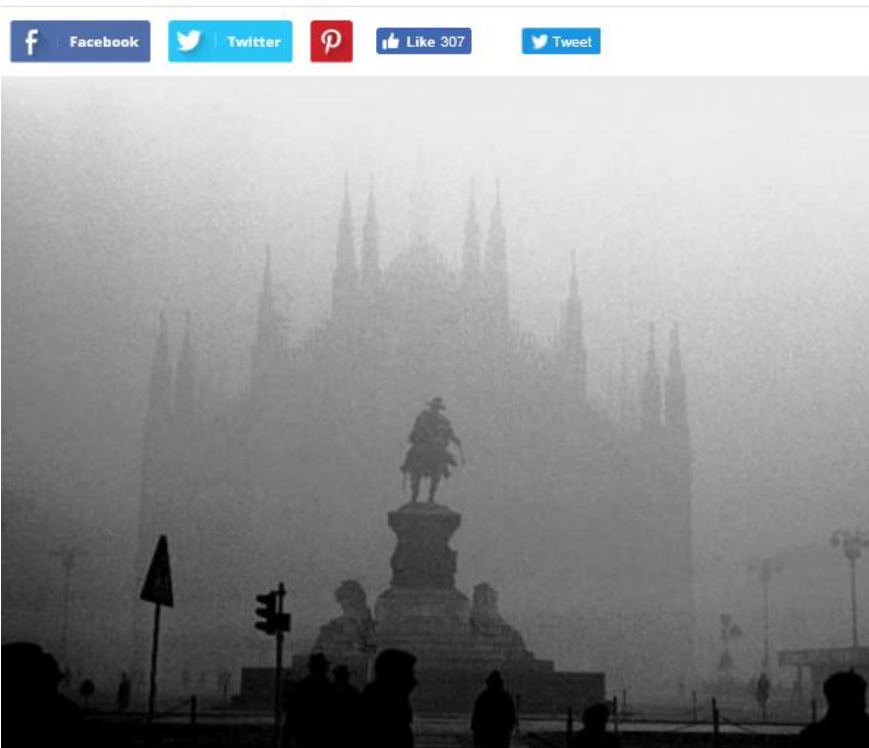


Milan has poor air quality

Motivation

Milan has second worst smog in Europe – WHO

30 Jan, 2018



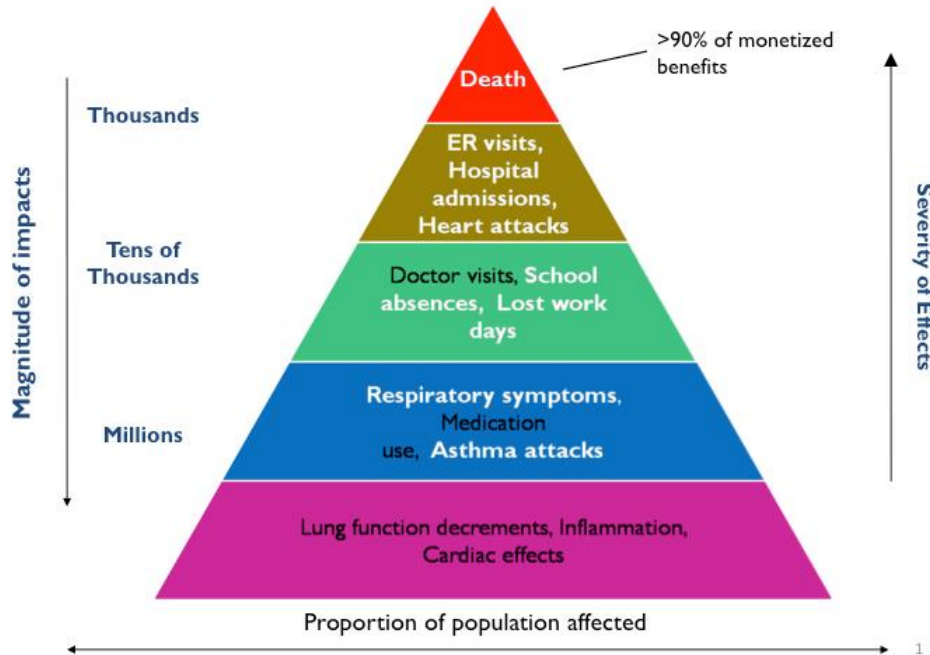
A report by the World Health Organization has placed Milan just behind Turin and just before Naples as the three European cities with the worst levels of atmospheric pollution.

Article based on WHO report using 2016 data

Poor air quality affects citizens' health

Motivation

A “Pyramid of Effects” from Air Pollution

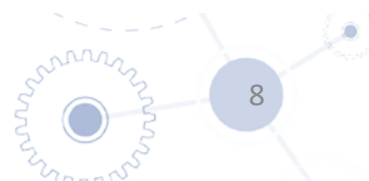


Fine particles can enter deep into the lungs and enter the blood stream. **Health impacts from particles include:**

- **Premature death**
- **Non-fatal heart attacks**
- **Aggravated asthma**

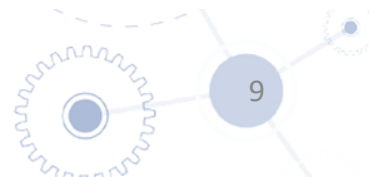
US Environmental Protection Agency

**Predicting AQI could help decision makers
introduce interventions that are predicted to
manage air quality in real time improving
citizens' health**



III. Data

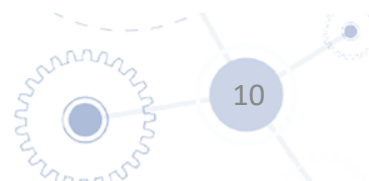
- I. Methodology
- II. Feature Construction
- III. Data Exploration



We will take sensor data from weather and traffic to predict AQI using machine learning methods, selecting on accuracy

Methodology

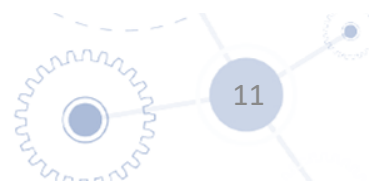
- The process will be the following:
 1. Features will be constructed, including the target
 2. Features will be selected based on a variety of dimensionality reduction techniques and intuition
 3. Train and test sets (Validation?) will be created
 4. Multiple prediction models will be trained and tested with these variables
 5. The model with the most accurate predictions in the test period will be selected



Data had to be extracted from files generated by sensors, cleaned and merged

Feature Construction

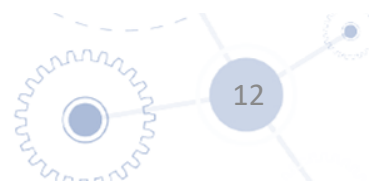
- The data had three main sources
 - Weather Sensors
 - Traffic Gates
 - Pollution Sensors (used to calculate AQI)
- Sensors measuring the same thing were averaged
- Source of this data was X



Missing data was an issue

Feature Construction

- We wished to predict on the hourly level of data, which was not always available. Additionally, in some periods data was just missing. The data was either recorded hourly for the whole period, or imputed to hourly with the appropriate covariates
- Noisy time series data with potential for measurement errors was smoothed ($\alpha=.2$)



Some features were constructed somewhat arbitrarily

Feature Construction

- The vehicle length was selected by visually guessing where a tri modal distribution was best separated into small, medium and large cars
- We took the 15 minute traffic counts






AQI had to be calculated and then transformed to a classifier to serve as the target variable

Feature Construction

- Everybody's favorite equation

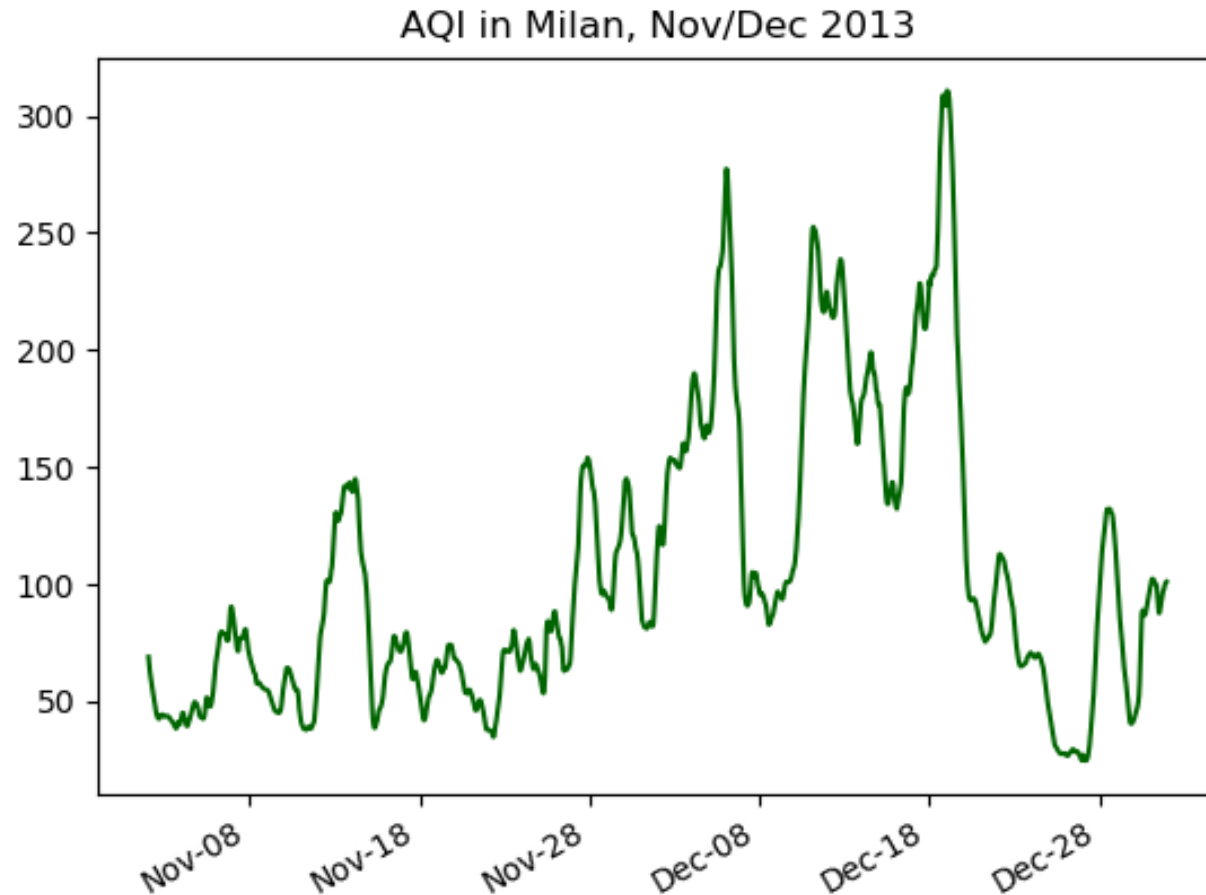
$$I_{QA} = \frac{I_{PM10} + \max(I_{NO2}, I_{O3})}{2}$$

- Graphic of Scale from poor to good

Valori dell'indice	Cromatismi	Qualità dell'aria
< 50		Buona
50-99		Accettabile
100-149		Mediocre
150-199		Scadente
> 200		Pessima

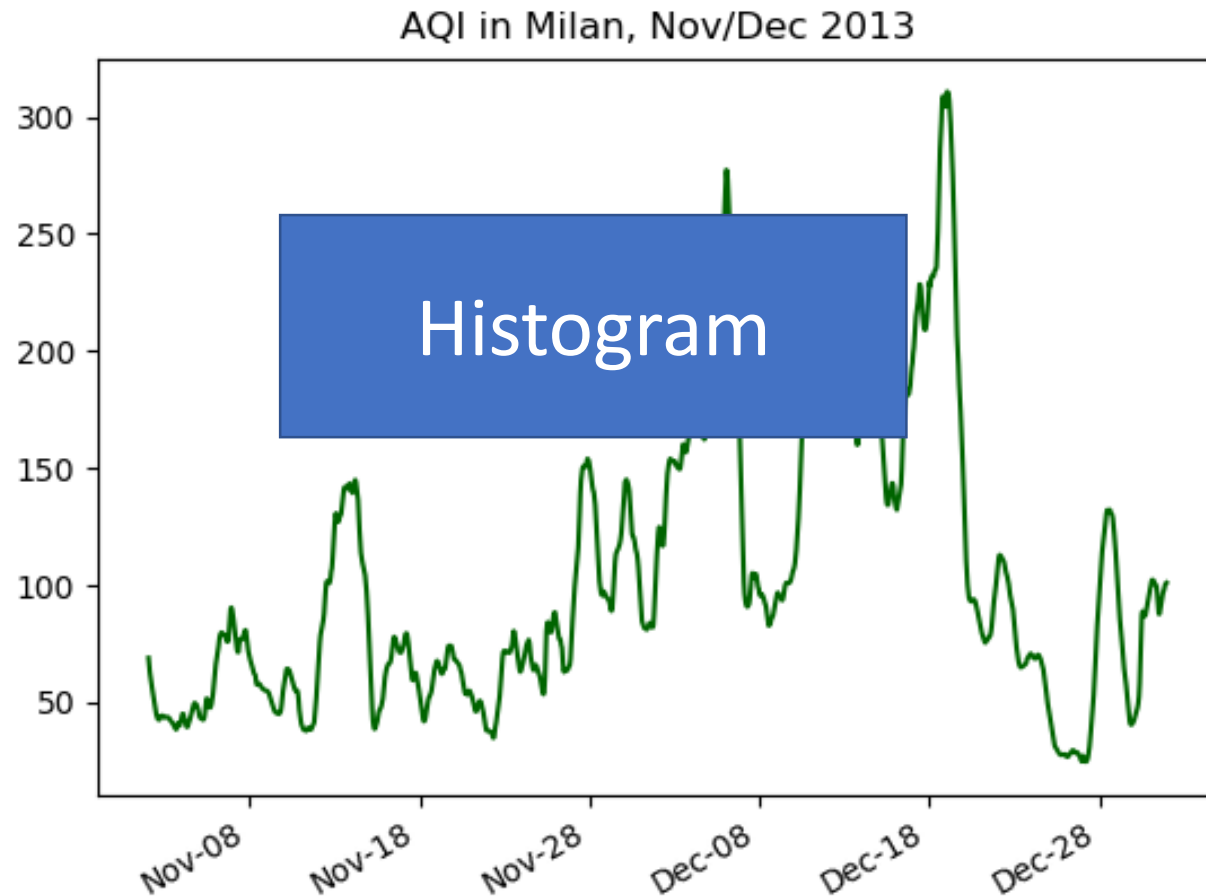
Air Quality Index score

Feature Construction



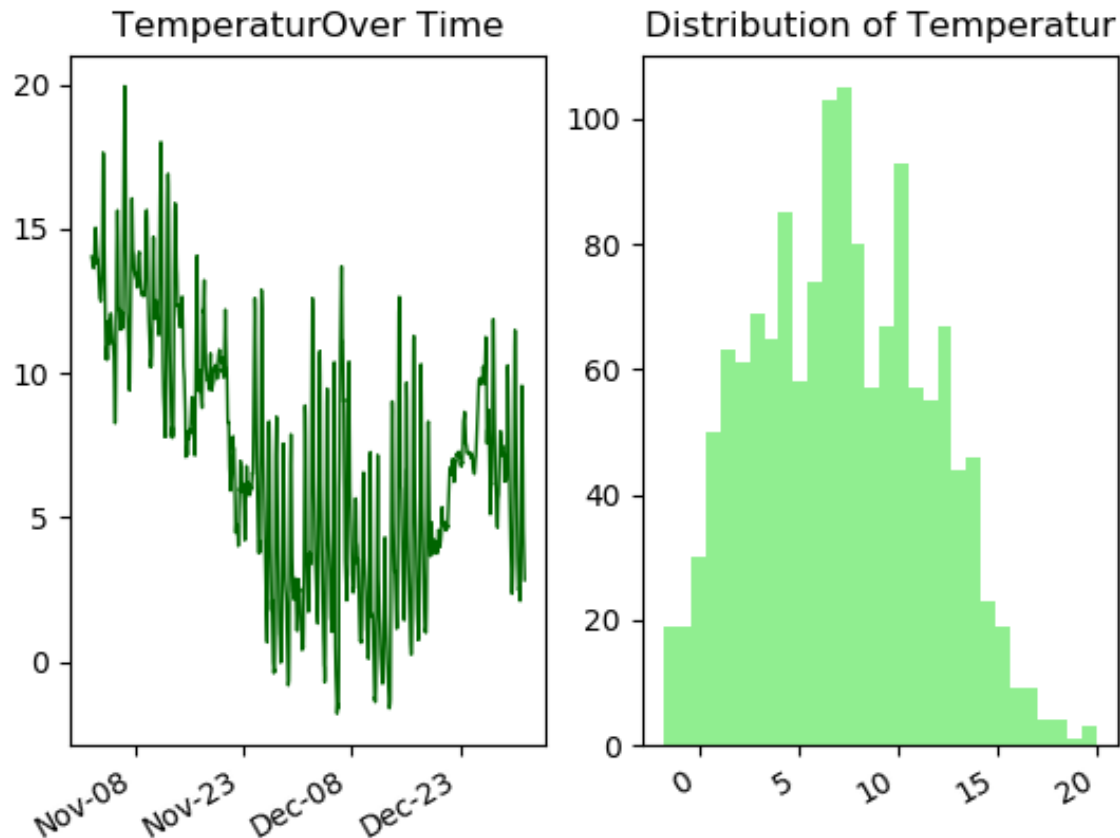
Air Quality Index score distribution

Feature Construction



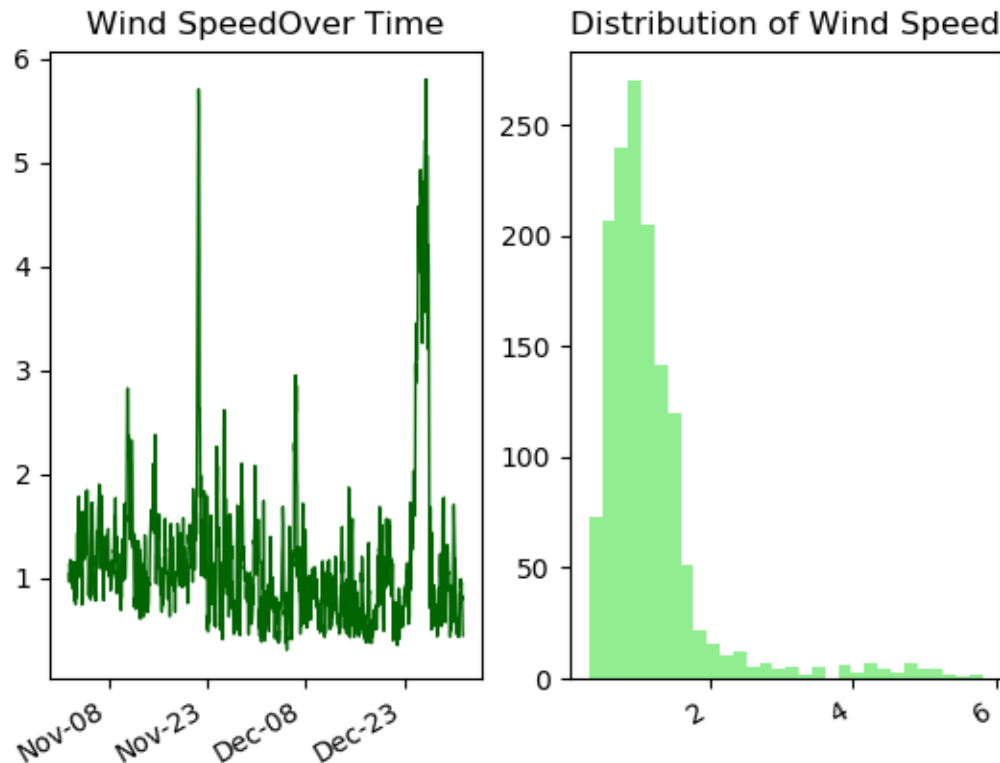
Temperature had a normal distribution and declined into year end

Data Exploration



Wind speed was very right skewed, with large peaks

Data Exploration



Traffic had obvious daily/weekly pattern and a “tri-modal” distribution

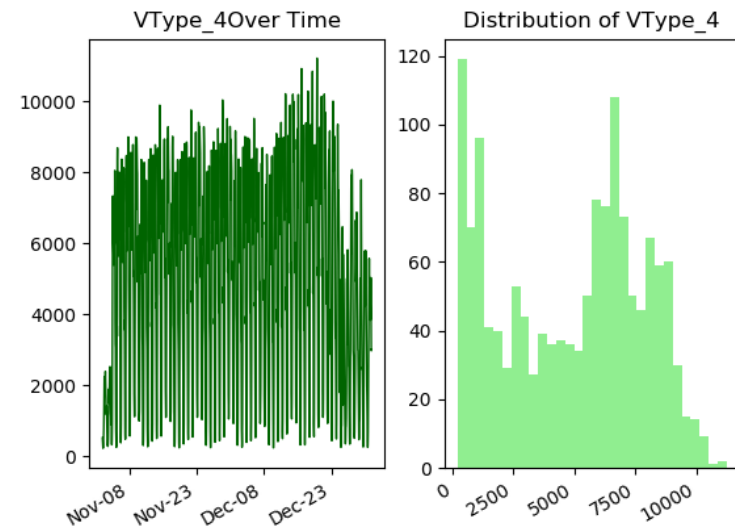
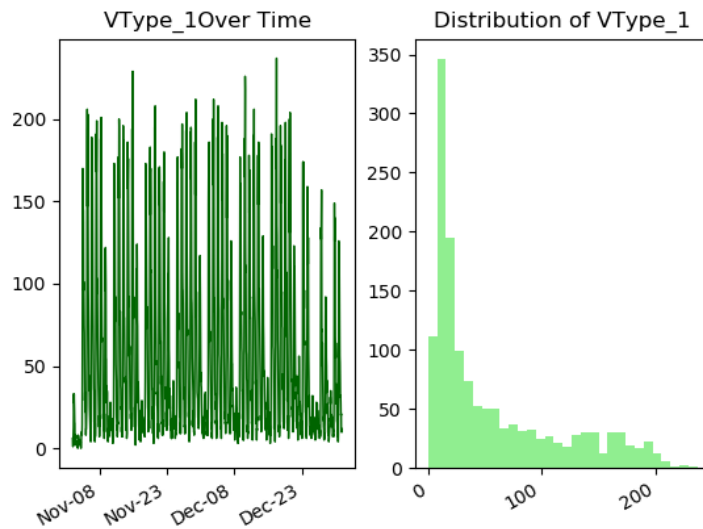
Data Exploration



Different Vehicle types had different distributions

Data Exploration

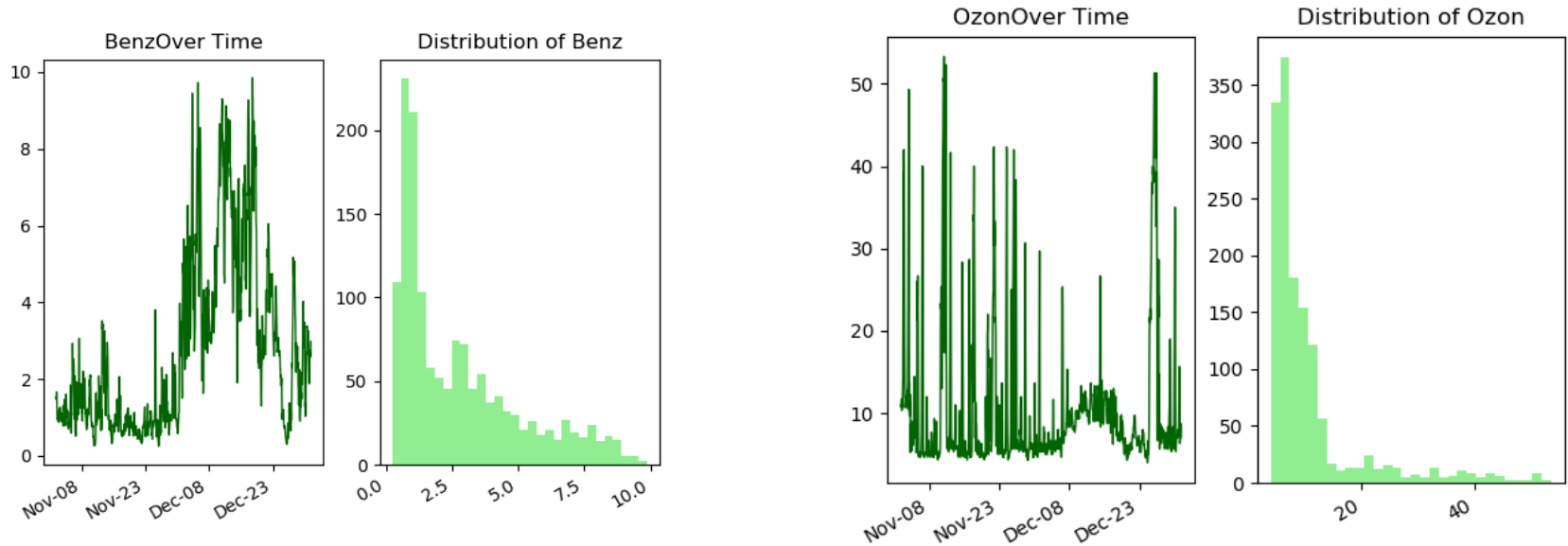
- Can see weekends in 1
- Dist more bimodal or uniform in 4, more power lawish in 1
- What are the vtypes and can we tell a story?



Pollutants looked a lot like AQI pollutants, Ozone least so

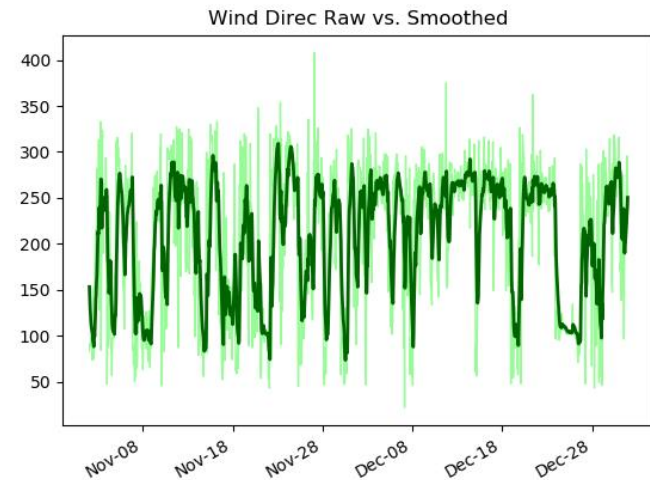
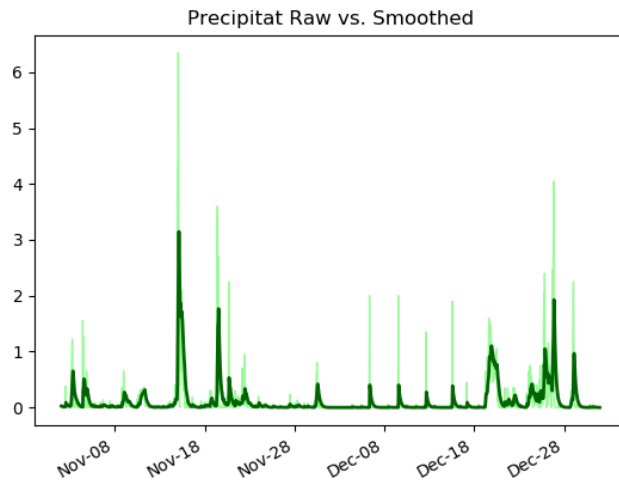
Data Exploration

- Show example of correlated, least correlated, most anticorrelated pollutants



Smoothing eliminates noise from the features

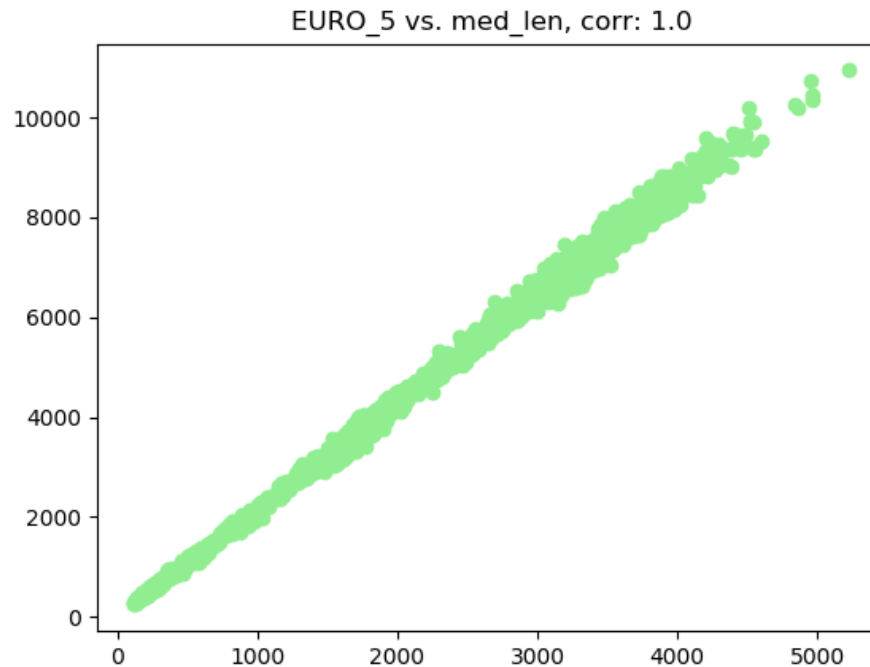
Data Exploration



Apparently Euro5 vehicles tend to be medium length

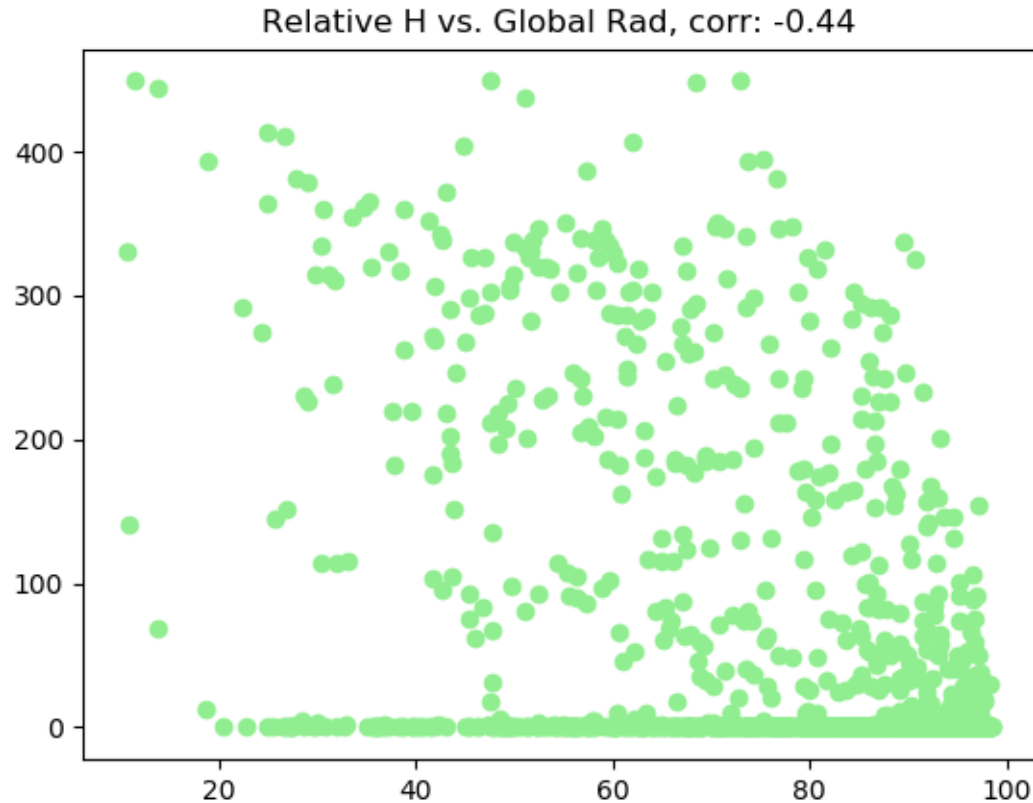
Data Exploration

- Though correlation rounds to 1, coefficient is about .5 so not 1:1
- Still, a good model should probably not have both of these variables



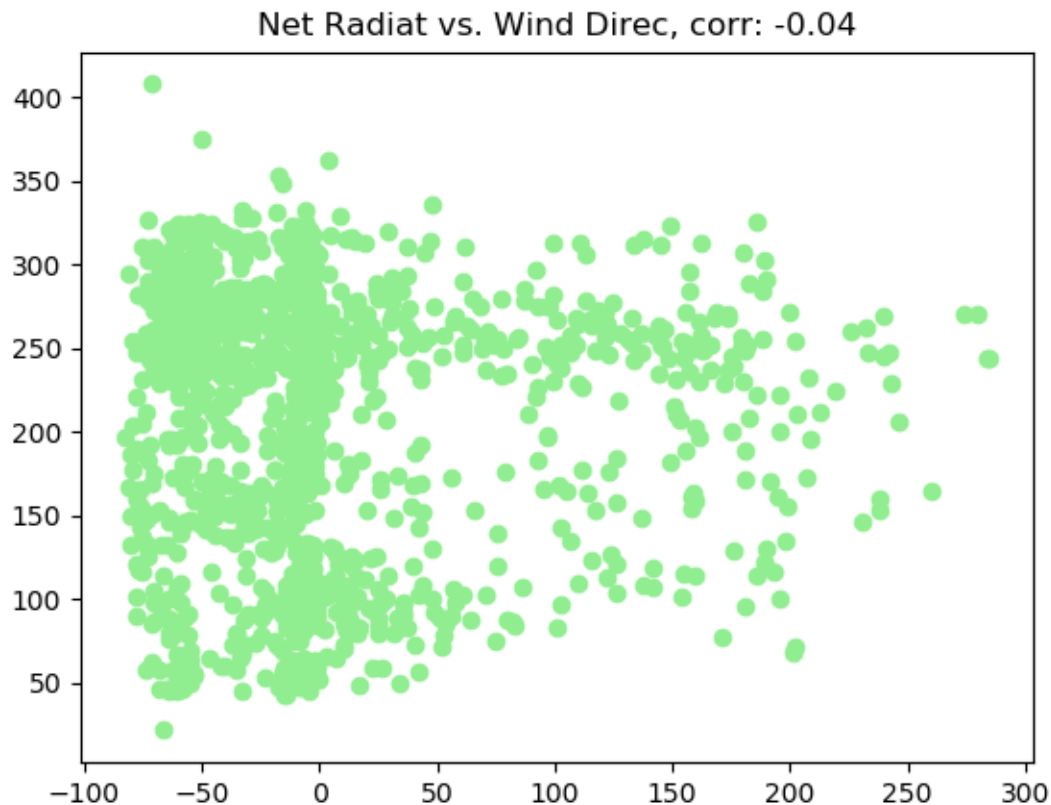
There were negatively correlated variables, but not as striking

Data Exploration



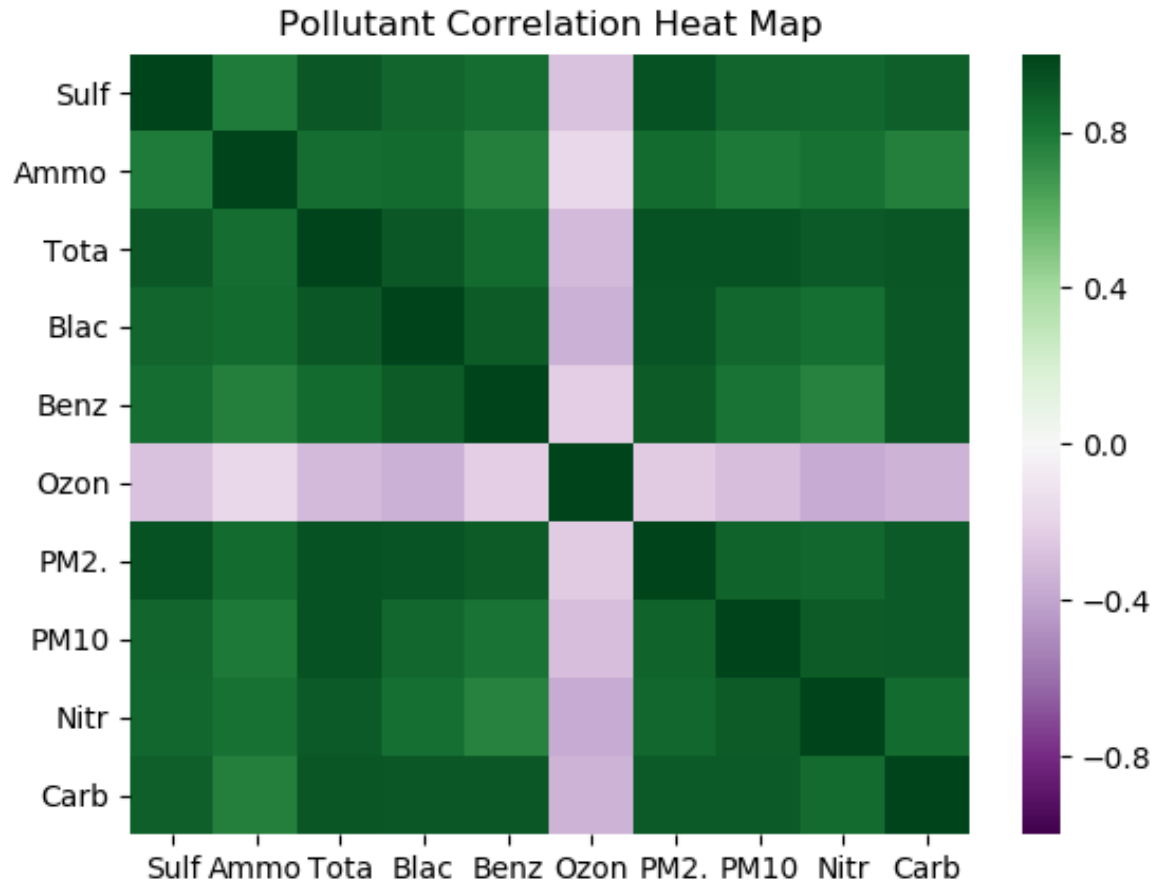
Some variables not correlated

Data Exploration



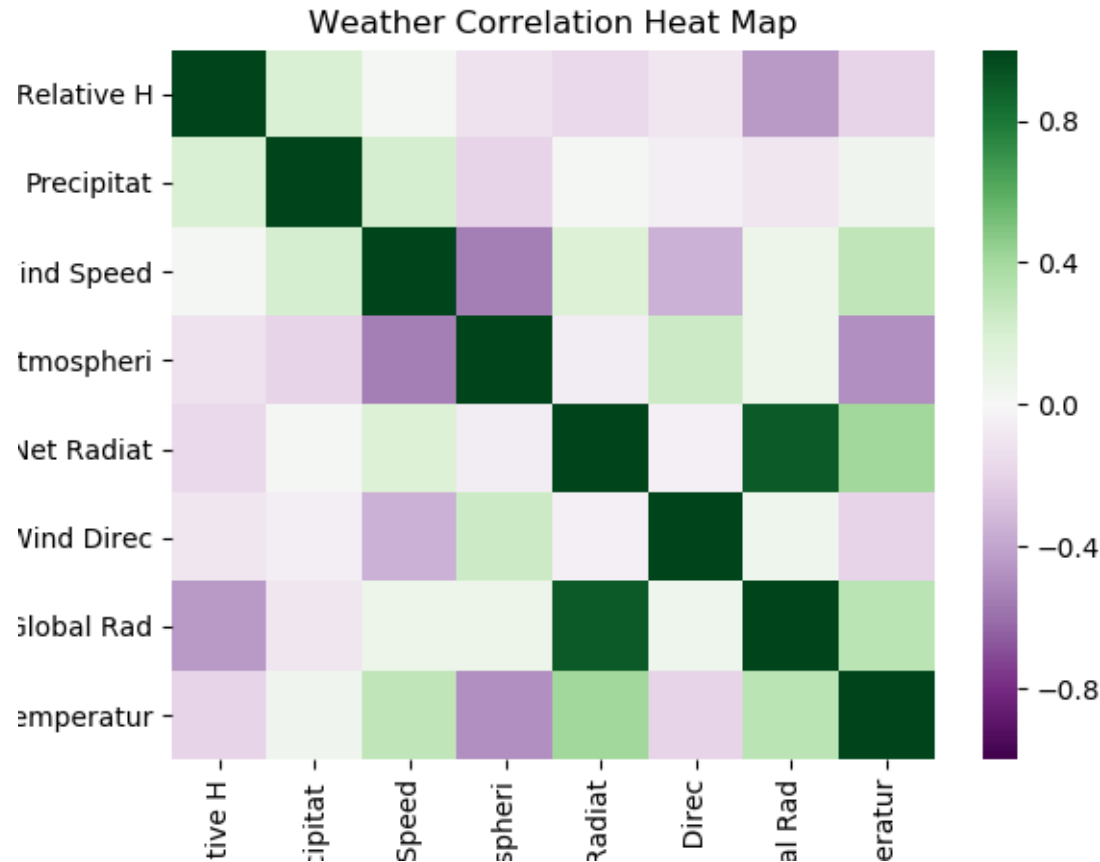
Pollutants are very correlated except ozone. This is good for imputation

Data Analysis - Multivariate



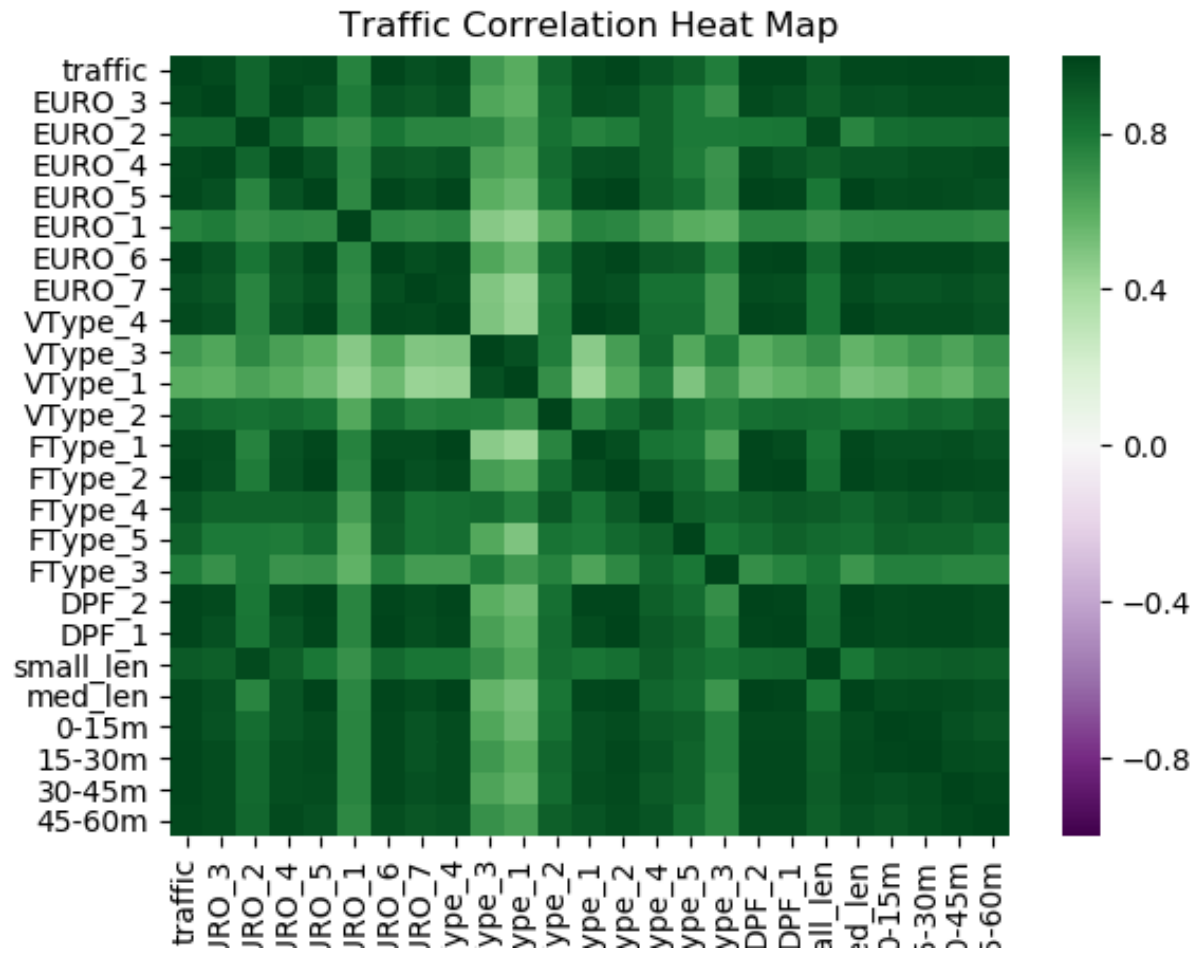
Weather data is not super correlated except the two radiation measures, which is to be expected

Data Exploration

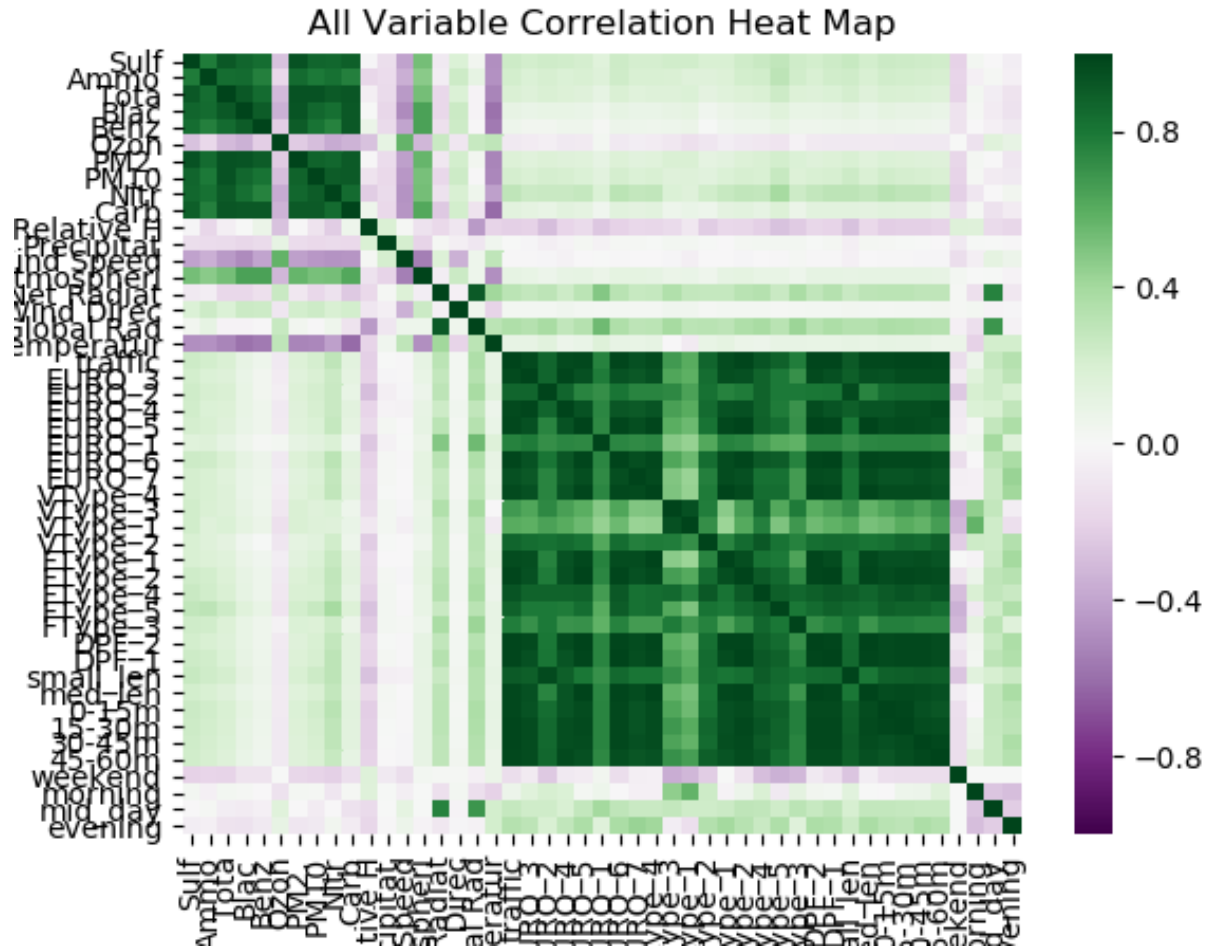


Traffic data is more correlated, and all positively correlated

Data Exploration

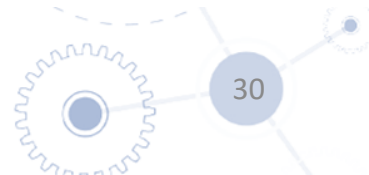


All data (should make this just features?)



IV. Model

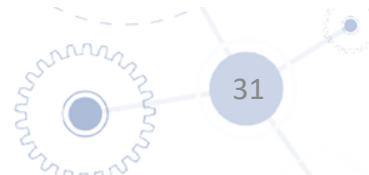
- I. Feature Selection
- II. Prediction
- III. Bonus: Is smoothed data better?



We have too many variables to drive an actionable result

Feature Selection

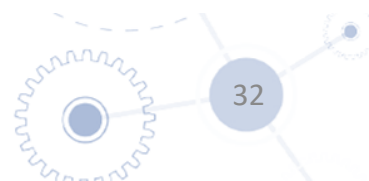
- Some techniques to do dimensionality reduction:
 - Penalization: Lasso
 - Embedded: Tree-based methods
 - Algorithmic: Stepwise selection
 - Heuristic: A method I made up
- After these we may apply a Voting: What variables are common to these methods?



Variables from each method

Feature Selection

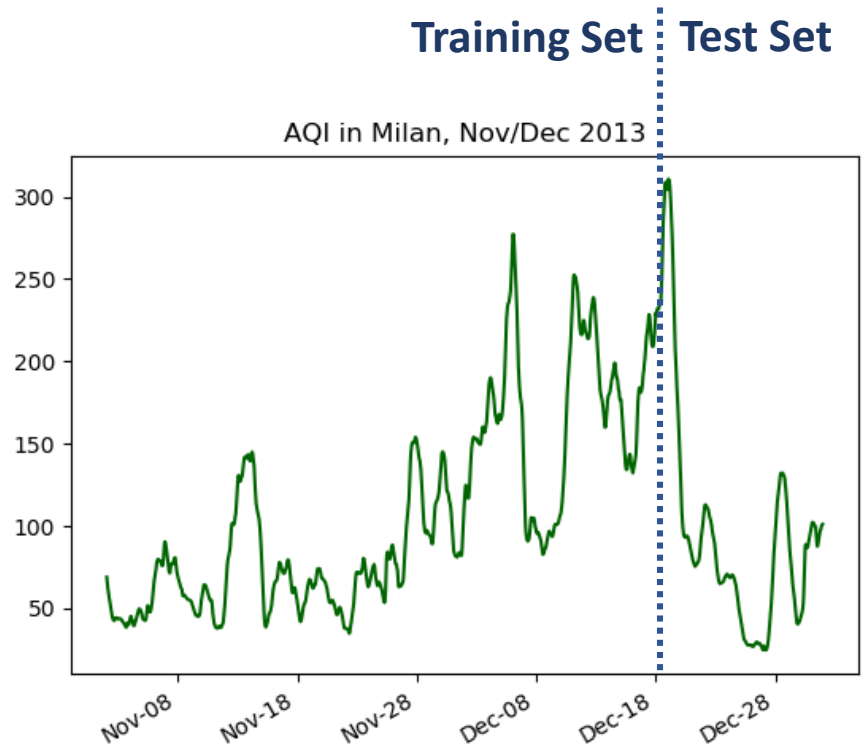
- Table
- Lets select the variables we want to use and try some methods that don't select for variables as well



Before we predict we need a train and test set

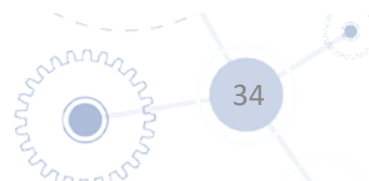
Prediction

- Distribution of classes in train test
- Strengths weaknesses tradeoff
- Validation?
- Visualization – histogram,
 - grey version of AQI to emphasize
 - dividing line?



Picking the modal training classification for the test data sets a “naive” benchmark

Prediction



Each of our variables alone does not match the “naive” accuracy

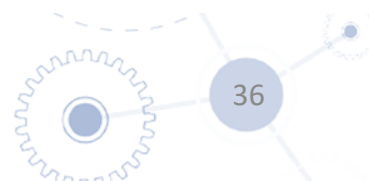
Prediction

- Bar chart of variables' accuracy

Each of our variables alone does not match the “naive” accuracy

Prediction

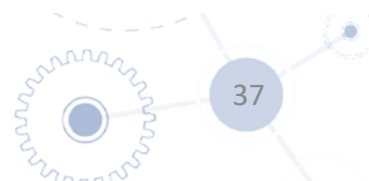
- Time series visualization of test data



Using the variables together yields better results, best in random forest probably with X%

Prediction

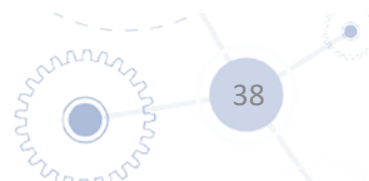
- asdf



Using the variables together yields better results

Prediction

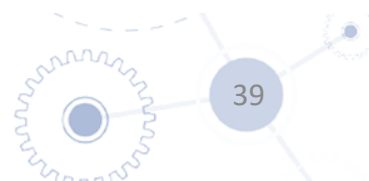
- Visualization of test set



Final selection of data and model

Prediction

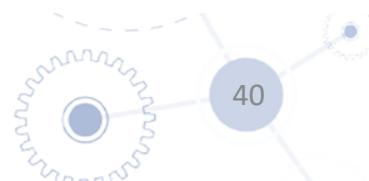
- asdf



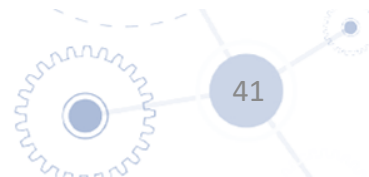
Rerunning the data on the nonsmoothed versions was better or worse

Noisy Data

- Wow I am surprised



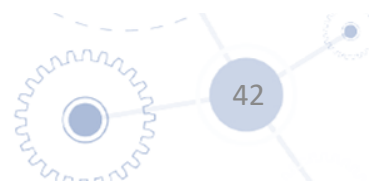
V. Conclusion



X model with Y variables performed prediction best

Conclusion

- Why it was selected
- Weaknesses
- Noisy data conclusion



Resources

