

Analysis of Air Quality in the United States

DS 5110 Spring 2020

Term Project

14 April, 2020



Carson Cook
Erin Keough
Paola Reyes
Aditya Samaroo
Priyank Shelat

Where does the data come from?

- Air Quality System (AQS) DataMart from the EPA
- Measurements of criteria pollutants
 - ▶ Ambient air pollution measurements from thousands of state and local monitoring locations
 - ▶ Particulate matter, sulfur dioxide, carbon monoxide, nitrogen dioxide, ozone, air toxics, lead
 - ▶ Hourly, 8-hour, daily, and annual concentrations
- Used by research, regulatory, and healthcare research communities
 - ▶ Clean Air Act compliance



<https://www.epa.gov>

Where does the data come from?



**CENTERS FOR DISEASE
CONTROL AND PREVENTION**

<https://www.cdc.gov>

- National Environmental Public Health Data tracking program from the CDC
- Data from a network of partners
 - ▶ National organizations, federal agencies, and fellowship participants
- Health and environmental data from city, state, and national sources
 - ▶ Air pollutant concentrations at county- and state-level resolutions
 - ▶ Prevalence of health conditions at county- and state-level resolutions
- Data available from 1999 - 2016

A look at the EPA data

state_code	county_code	site_num	parameter_code	pcu	latitude	longitude
Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705
Class :character	Class :character	Class :character	Class :character	1st Qu.: 1.000	Median : 1.000	Mean : 1.000
Mode :character	Mode :character	Mode :character	Mode :character	1st Qu.: 1.000	Median : 1.000	Mean : 1.000
Median : 1.000	Median : 1.000	Median : 1.000	Median : 1.000	Median : 1.000	Median : 1.000	Median : 1.000
1st Qu.: 1.000	1st Qu.: 1.000	1st Qu.: 1.000	1st Qu.: 1.000	1st Qu.: 1.000	1st Qu.: 1.000	1st Qu.: 1.000
3rd Qu.: 5.000	3rd Qu.: 5.000	3rd Qu.: 5.000	3rd Qu.: 5.000	3rd Qu.: 5.000	3rd Qu.: 5.000	3rd Qu.: 5.000
Max. : 168.64	Max. : 168.64	Max. : 168.64	Max. : 168.64	Max. : 168.64	Max. : 168.64	Max. : 168.64
Length:1624705	parameter_name	sample_duration	pollutant_code	metric_used	method_name	year
Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705
Class :character	Class :character	Class :character	Class :character	Class :character	Class :character	Class :character
Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character
Median : 2009	Median : 2009	Median : 2009	Median : 2009	Median : 2009	Median : 2009	Median : 2009
1st Qu.: 2009	1st Qu.: 2009	1st Qu.: 2009	1st Qu.: 2009	1st Qu.: 2009	1st Qu.: 2009	1st Qu.: 2009
3rd Qu.: 2014	3rd Qu.: 2014	3rd Qu.: 2014	3rd Qu.: 2014	3rd Qu.: 2014	3rd Qu.: 2014	3rd Qu.: 2014
Max. : 2019	Max. : 2019	Max. : 2019	Max. : 2019	Max. : 2019	Max. : 2019	Max. : 2019
units_of_measure	event_type	observation_count	observation_percent	completeness_indicator	valid_day_count	required_day_count
Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705
Class :character	Class :character	Class :character	Class :character	Class :character	Class :character	Class :character
Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character
Median : 95.00	Median : 95.00	Median : 95.00	Median : 95.00	Median : 60.00	Median : 60.00	Median : 60.00
Mean : 167.7	Mean : 82.74	Mean : 100.00	Mean : 100.00	Mean : 118.9	Mean : 145.0	Mean : 145.0
1st Qu.: 59.00	1st Qu.: 59.00	1st Qu.: 59.00	1st Qu.: 59.00	1st Qu.: 57.00	1st Qu.: 57.00	1st Qu.: 57.00
3rd Qu.: 100.00	3rd Qu.: 100.00	3rd Qu.: 100.00	3rd Qu.: 100.00	3rd Qu.: 100.00	3rd Qu.: 100.00	3rd Qu.: 100.00
Max. : 105025	Max. : 105025	Max. : 105025	Max. : 105025	Max. : 366.0	Max. : 366.0	Max. : 366.0
exceptional_data_count	null_data_count	primary_exceedance_count	secondary_exceedance_count	certification_indicator	num_obs_below_md	Length:1624705
Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705
Class :character	Class :character	Class :character	Class :character	Class :character	Class :character	Class :character
Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character
Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00
Mean : 30.89	Mean : 69.48	Mean : 2.6	Mean : 2.6	Mean : 2.6	Mean : 0	Mean : 0
1st Qu.: 0.00	1st Qu.: 0.00	1st Qu.: 0.00	1st Qu.: 0.00	1st Qu.: 0.00	1st Qu.: 0.00	1st Qu.: 0.00
3rd Qu.: 102471.00	3rd Qu.: 105096.00	3rd Qu.: 372.00	3rd Qu.: 372.00	3rd Qu.: 0.00	Max. : 105096.00	Max. : 105096.00
Max. : 105096.00	Max. : 105096.00	Max. : 105096.00	Max. : 105096.00	Max. : 0.00	NA's : 1375186	NA's : 1394666
arithmetic_mean	arithmetic_standard_dev	first_max_value	second_max_value	second_max_datetime	third_max_value	third_max_datetime
Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705
Class :character	Class :character	Class :character	Class :character	Class :character	Class :character	Class :character
Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character
Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00
Mean : 30.89	Mean : 69.48	Mean : 2.6	Mean : 2.6	Mean : 2.6	Mean : 0	Mean : 0
1st Qu.: 0.00	1st Qu.: 0.00	1st Qu.: 0.00	1st Qu.: 0.00	1st Qu.: 0.00	1st Qu.: 0.00	1st Qu.: 0.00
3rd Qu.: 102471.00	3rd Qu.: 105096.00	3rd Qu.: 372.00	3rd Qu.: 372.00	3rd Qu.: 0.00	Max. : 105096.00	Max. : 105096.00
Max. : 105096.00	Max. : 105096.00	Max. : 105096.00	Max. : 105096.00	Max. : 0.00	NA's : 130496	NA's : 140496
first_no_max_datetime	second_max_non_overlapping_value	fourth_max_datetime	fourth_max_value	fourth_max_datetime	fourth_max_non_overlapping_value	third_max_value
Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705
Class :character	Class :character	Class :character	Class :character	Class :character	Class :character	Class :character
Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character
Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00
Mean : 24.20	Mean : 24.20	Mean : 24.20	Mean : 24.20	Mean : 24.20	Mean : 24.20	Mean : 24.20
1st Qu.: 0.07	1st Qu.: 0.07	1st Qu.: 0.07	1st Qu.: 0.07	1st Qu.: 0.07	1st Qu.: 0.07	1st Qu.: 0.07
3rd Qu.: 24.40	3rd Qu.: 24.40	3rd Qu.: 24.40	3rd Qu.: 24.40	3rd Qu.: 24.40	3rd Qu.: 24.40	3rd Qu.: 24.40
Max. : 24996.00	Max. : 24996.00	Max. : 24996.00	Max. : 24996.00	Max. : 24996.00	Max. : 24996.00	Max. : 24996.00
NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496
first_no_max_datetime	second_max_non_overlapping_value	second_no_max_datetime	second_no_max_value	second_no_max_datetime	second_no_max_value	ninety_nine_percentile_ninety_eight_percentile
Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705
Class :character	Class :character	Class :character	Class :character	Class :character	Class :character	Class :character
Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character
Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00	Median : 0.00
Mean : 2003-11-30 08:45:00	Mean : 2003-11-30 08:45:00	Mean : 2003-11-30 08:45:00	Mean : 2003-11-30 08:45:00	Mean : 2003-11-30 08:45:00	Mean : 2003-11-30 08:45:00	Mean : 2003-11-30 08:45:00
1st Qu.: 2008-01-12 00:30:00	1st Qu.: 2008-01-12 00:30:00	1st Qu.: 2008-01-12 00:30:00	1st Qu.: 2008-01-12 00:30:00	1st Qu.: 2008-01-12 00:30:00	1st Qu.: 2008-01-12 00:30:00	1st Qu.: 2008-01-12 00:30:00
3rd Qu.: 2013-11-29 06:00:00	3rd Qu.: 2013-11-29 06:00:00	3rd Qu.: 2013-11-29 06:00:00	3rd Qu.: 2013-11-29 06:00:00	3rd Qu.: 2013-11-29 06:00:00	3rd Qu.: 2013-11-29 06:00:00	3rd Qu.: 2013-11-29 06:00:00
Max. : 2019-10-26 01:00:00	Max. : 2019-10-26 01:00:00	Max. : 2019-10-26 01:00:00	Max. : 2019-10-26 01:00:00	Max. : 2019-10-26 01:00:00	Max. : 2019-10-26 01:00:00	Max. : 2019-10-26 01:00:00
NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496
ninetyFive_percentile_ninetyPercentile	seventyFive_percentile_fiftyPercentile	tenPercentile	local_site_name	date_of_last_change	Length:1624705	Length:1624705
Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705
Class :character	Class :character	Class :character	Class :character	Class :character	Class :character	Class :character
Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character
Median : 1.02	Median : 0.80	Median : 0.50	Median : 0.50	Median : 0.50	Median : 0.50	Median : 0.50
Mean : 63.00	Mean : 39.00	Mean : 24.00	Mean : 24.00	Mean : 24.00	Mean : 24.00	Mean : 24.00
1st Qu.: 19.60	1st Qu.: 16.00	1st Qu.: 11.10	1st Qu.: 11.10	1st Qu.: 11.10	1st Qu.: 11.10	1st Qu.: 11.10
3rd Qu.: 45.10	3rd Qu.: 34.00	3rd Qu.: 24.00	3rd Qu.: 24.00	3rd Qu.: 24.00	3rd Qu.: 24.00	3rd Qu.: 24.00
Max. : 110000.00	Max. : 110000.00	Max. : 110000.00	Max. : 110000.00	Max. : 110000.00	Max. : 110000.00	Max. : 110000.00
NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496
address	state_name	county_name	city_name	cbsa_name	date_of_last_change	Length:1624705
Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705	Length:1624705
Class :character	Class :character	Class :character	Class :character	Class :character	Class :character	Class :character
Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character	Mode :character
Median : 1.02	Median : 0.80	Median : 0.50	Median : 0.50	Median : 0.50	Median : 0.50	Median : 0.50
Mean : 63.00	Mean : 39.00	Mean : 24.00	Mean : 24.00	Mean : 24.00	Mean : 24.00	Mean : 24.00
1st Qu.: 19.60	1st Qu.: 16.00	1st Qu.: 11.10	1st Qu.: 11.10	1st Qu.: 11.10	1st Qu.: 11.10	1st Qu.: 11.10
3rd Qu.: 45.10	3rd Qu.: 34.00	3rd Qu.: 24.00	3rd Qu.: 24.00	3rd Qu.: 24.00	3rd Qu.: 24.00	3rd Qu.: 24.00
Max. : 110000.00	Max. : 110000.00	Max. : 110000.00	Max. : 110000.00	Max. : 110000.00	Max. : 110000.00	Max. : 110000.00
NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496	NA's : 140496

- 1.6 million rows of annual summary data
- Data can be roughly broken into:
 - Station information
 - Parameter information
 - Observation information

How is the EPA data collected?

- 5624 air monitoring locations in the contiguous US
- 1052 distinct parameters measured and recorded
- Air monitoring locations can measure one or more parameters



<https://www.mdaqmd.ca.gov/permitting/compliance>

Observations from raw EPA data

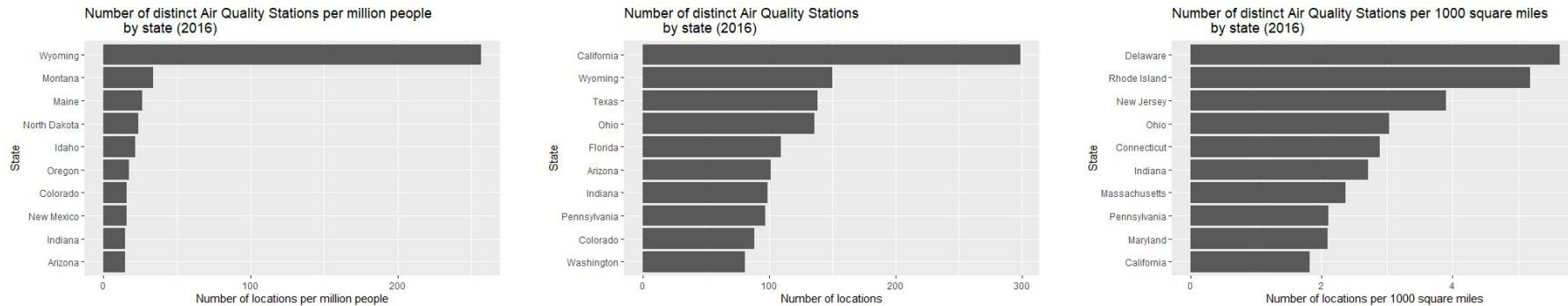
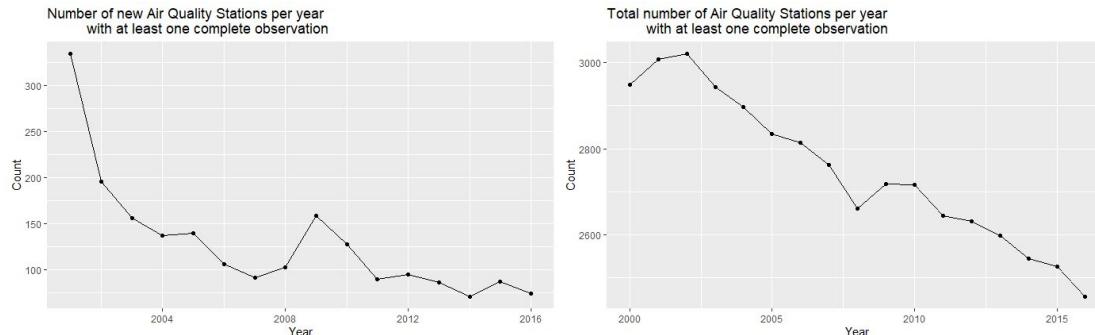


Figure 1. Number of Air Quality Stations per state does not seem to be dependent on population nor square mileage

Observations from raw EPA data

Any Air Quality Station



PM2.5 Air Quality Stations

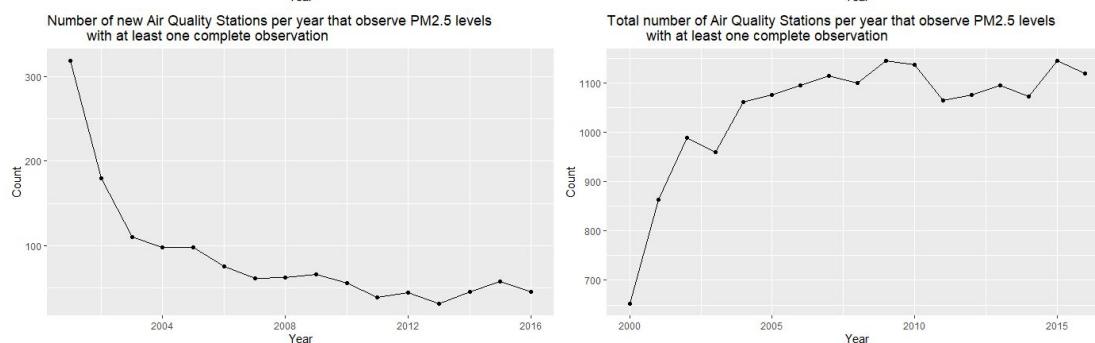


Figure 2. Initial analysis of the AQS dataset from the EPA.

Data Pre-Processing

Selection

- Filtered out variables with more than 50% missing values.
- Removed duplicates.

Preprocessing

- Tidy dataset: pivoted variables.
- Merged datasets: EPA + CDC
- Partition Dataset: 80% training 20% testing

Transformation

- Imputed missing values.
- Log-transformed both target and explanatory variables.

Exploratory Data Analysis - PM_{2.5}

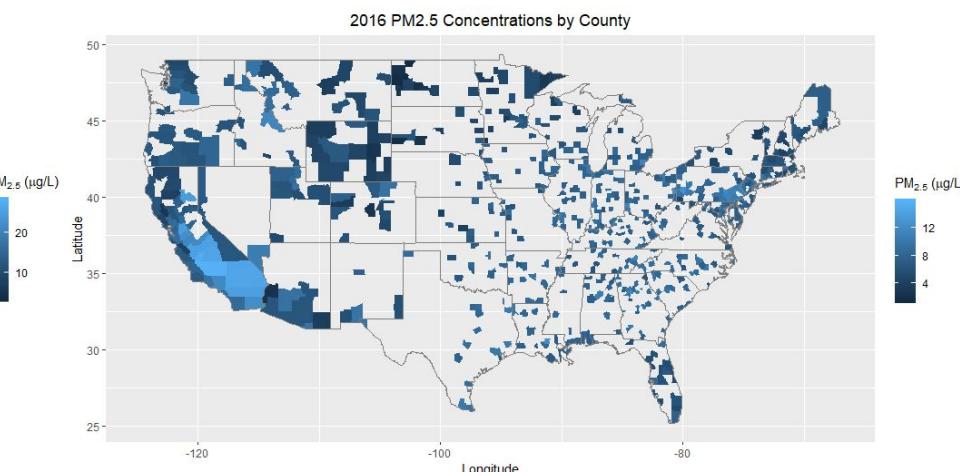
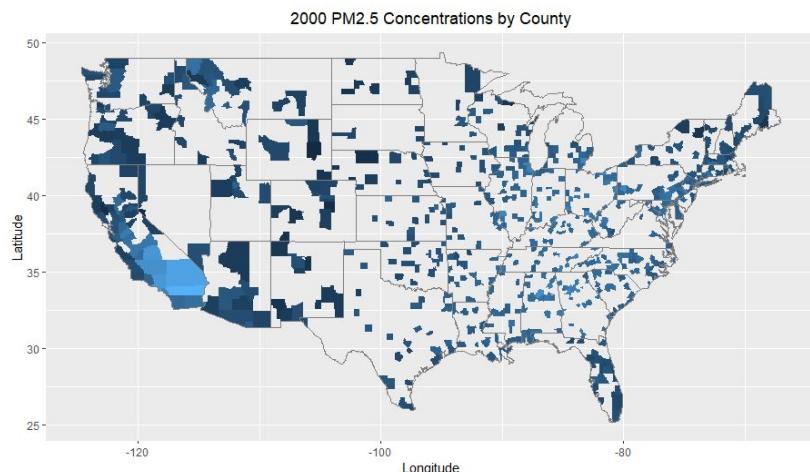


Figure 3. Comparison of particulate matter concentrations by county in 2000 and 2016.

Exploratory Data Analysis - PM_{2.5}

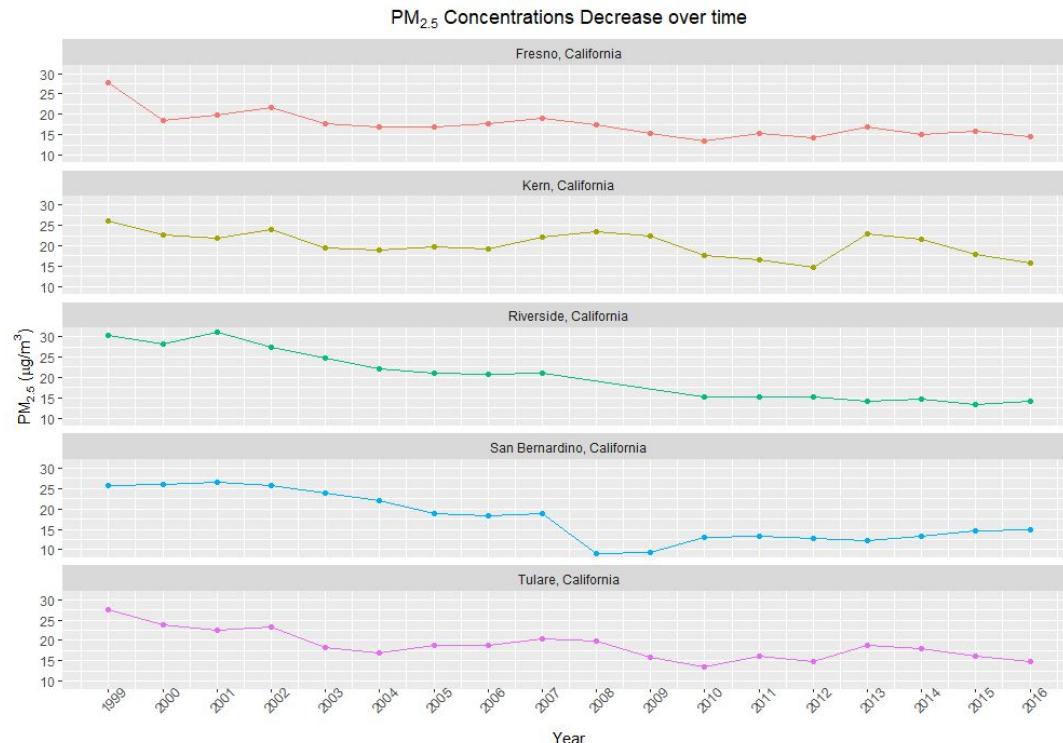


Table 1. Average, minimum, and maximum concentrations for each of the top 5 counties.

CountyID	Average	Maximum	Minimum
Riverside, California	20.53750	31.0	13.4
Kern, California	20.41667	26.2	14.7
Tulare, California	18.79444	27.6	13.6
San Bernardino, California	17.67222	26.5	9.1
Fresno, California	17.43333	27.7	13.4

Figure 4. PM_{2.5} concentrations in the areas with the highest overall concentrations from 1999 to 2016.

Exploratory Data Analysis - PM_{2.5} and Ozone

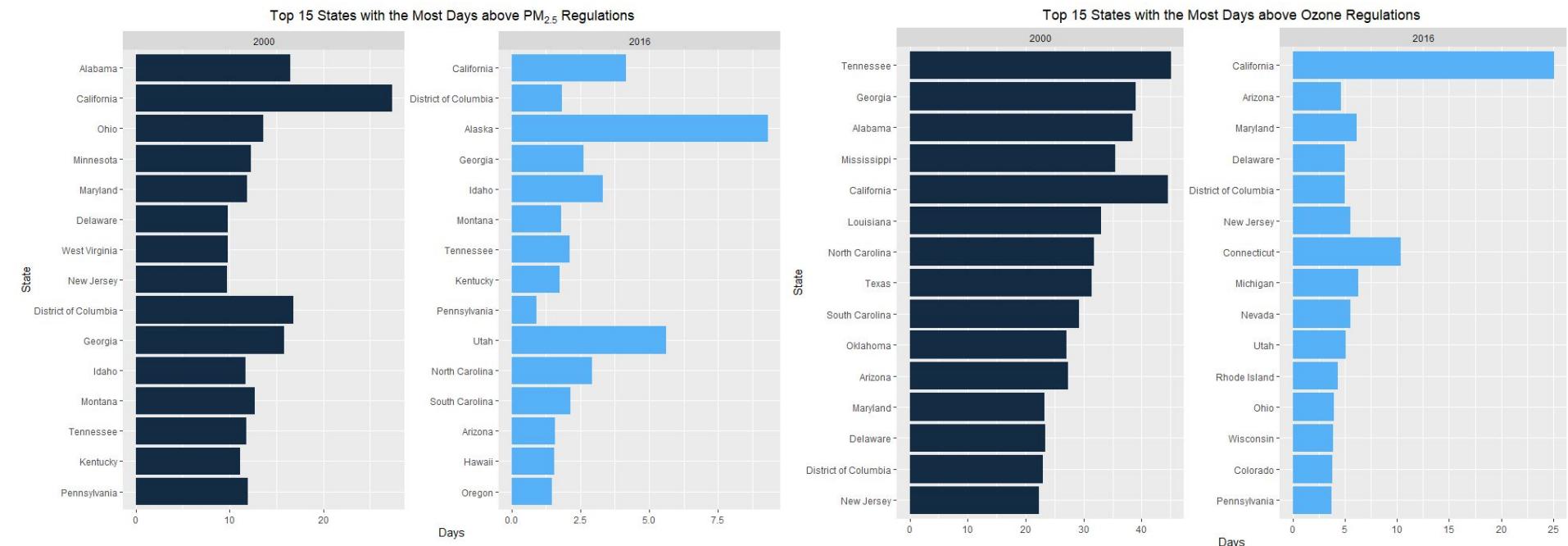


Figure 5. States with the most days above PM_{2.5} and Ozone regulation in 2000 and 2016.

Exploratory Data Analysis - PM_{2.5}

Reduction Percent of PM 2.5 <chr>	Minimum Mortality Benefit <dbl>	Maximum Mortality Benefit <dbl>	Average Mortality Benefit <dbl>	Median Mortality Benefit <dbl>
Reduction Percent: 10% Reduction	0	679	4.785960	1
Reduction Percent: 15% Reduction	0	958	7.155908	2
Reduction Percent: 20% Reduction	0	972	9.435453	3
Reduction Percent: 25% Reduction	0	942	11.599772	4
Reduction Percent: 5% Reduction	0	341	2.366039	1

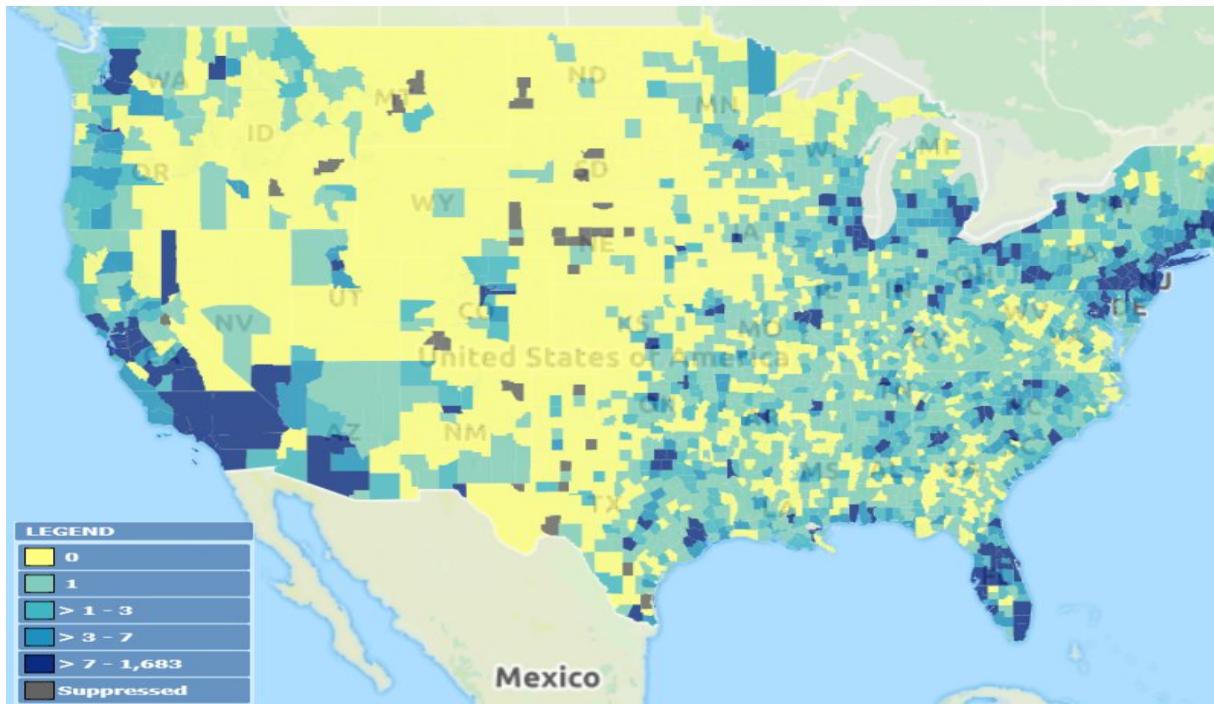


Table 2. PM_{2.5} Mortality Benefit Table.

Deaths Avoided by 5%
Reduction in PM 2.5



Top 5 Most Affected States

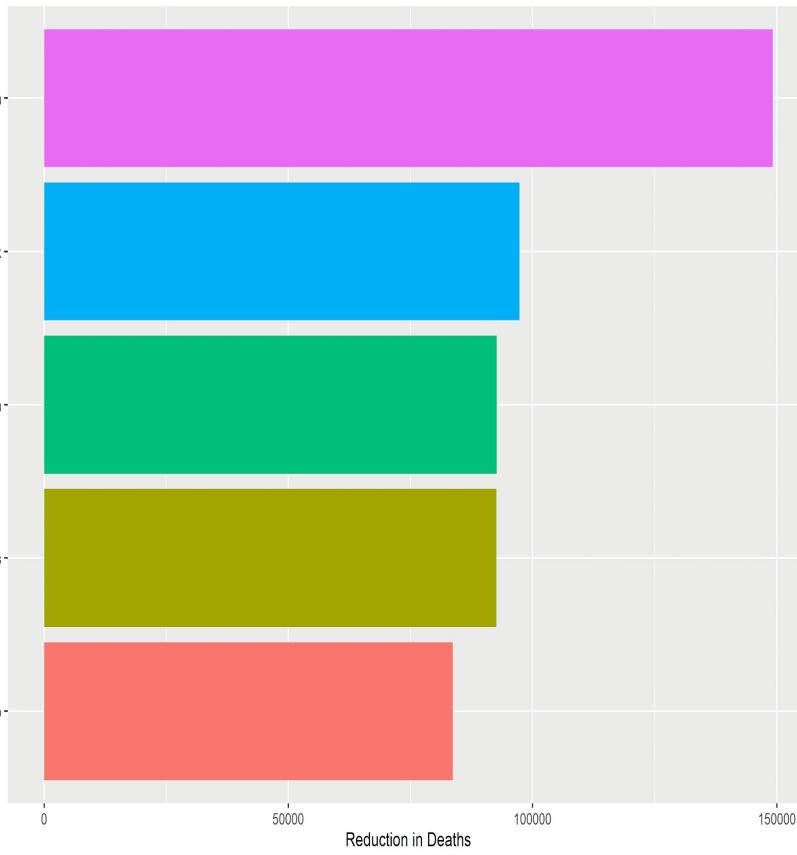


Figure 7. Top states impacted by reducing PM_{2.5} concentrations.

Mortality Benefit by State

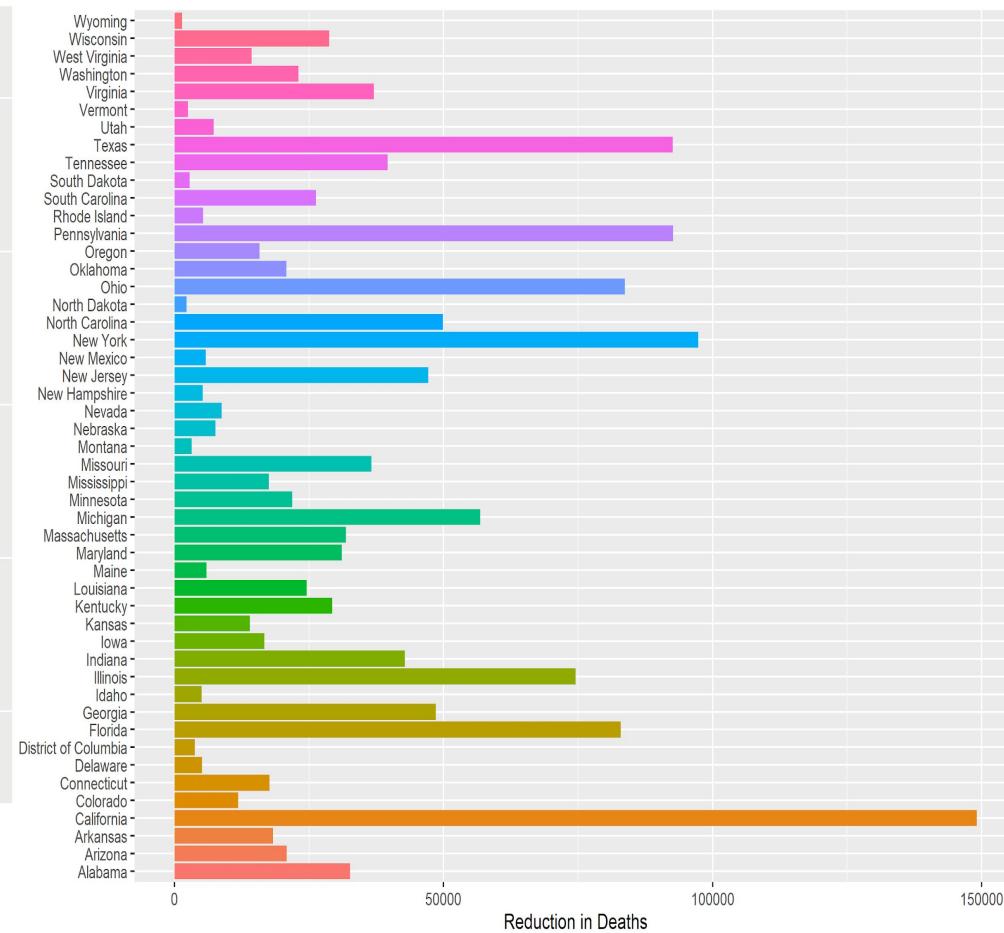


Figure 8. Benefits of reducing PM_{2.5} concentrations by state.

The Benefits of PM 2.5 are Decreasing Over Time

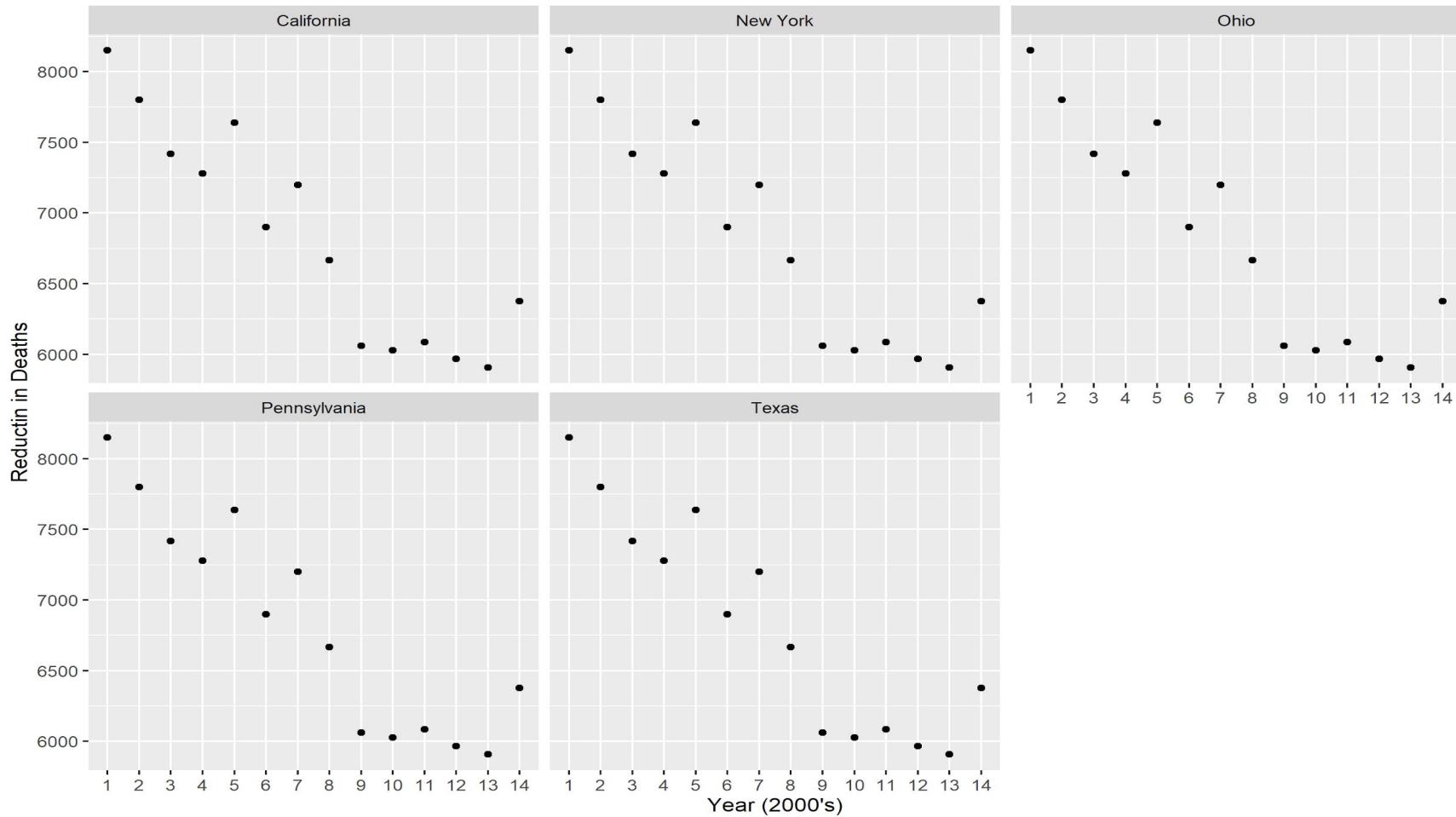


Figure 9. Mortality benefits of decreasing PM_{2.5} over time.

Exploratory Data Analysis - Asthma

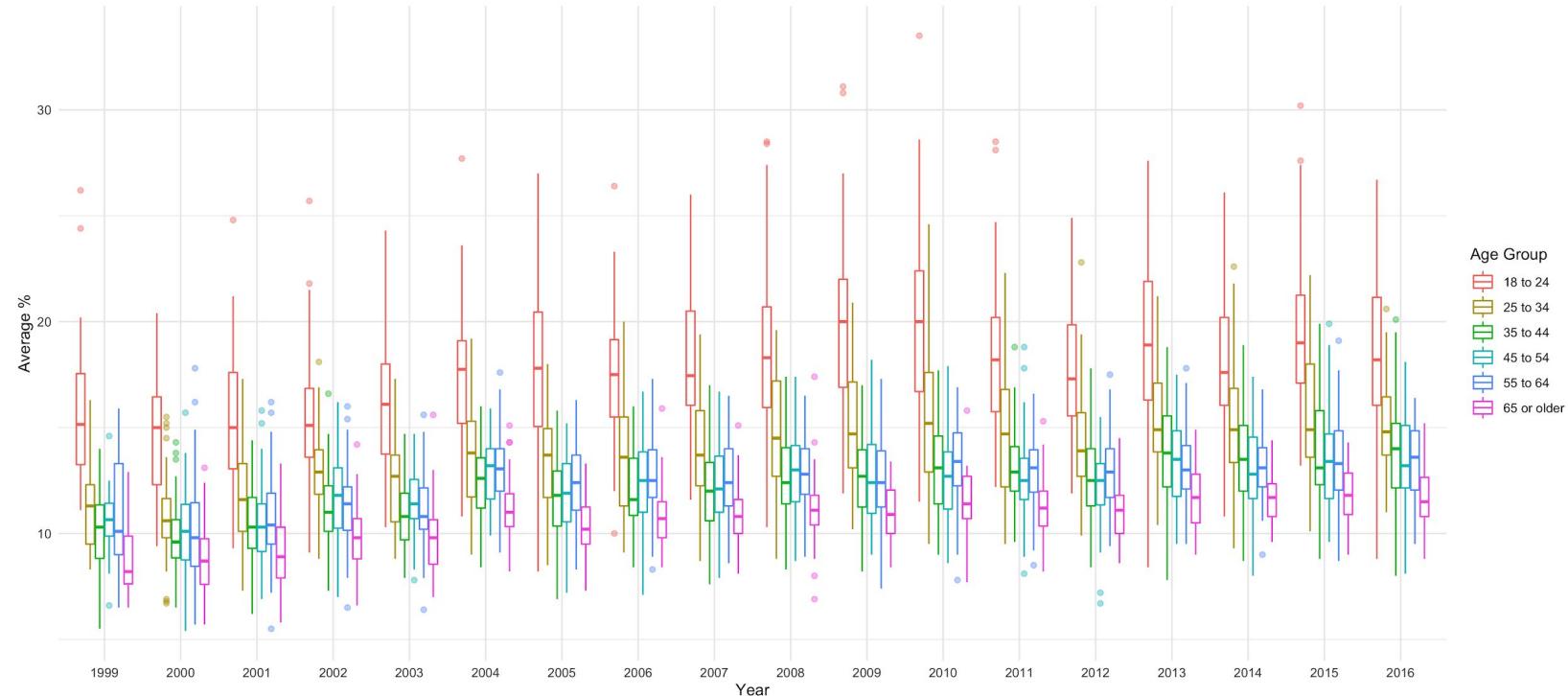


Figure 10. Asthma Prevalence by Age Group and Year.

Exploratory Data Analysis - Asthma

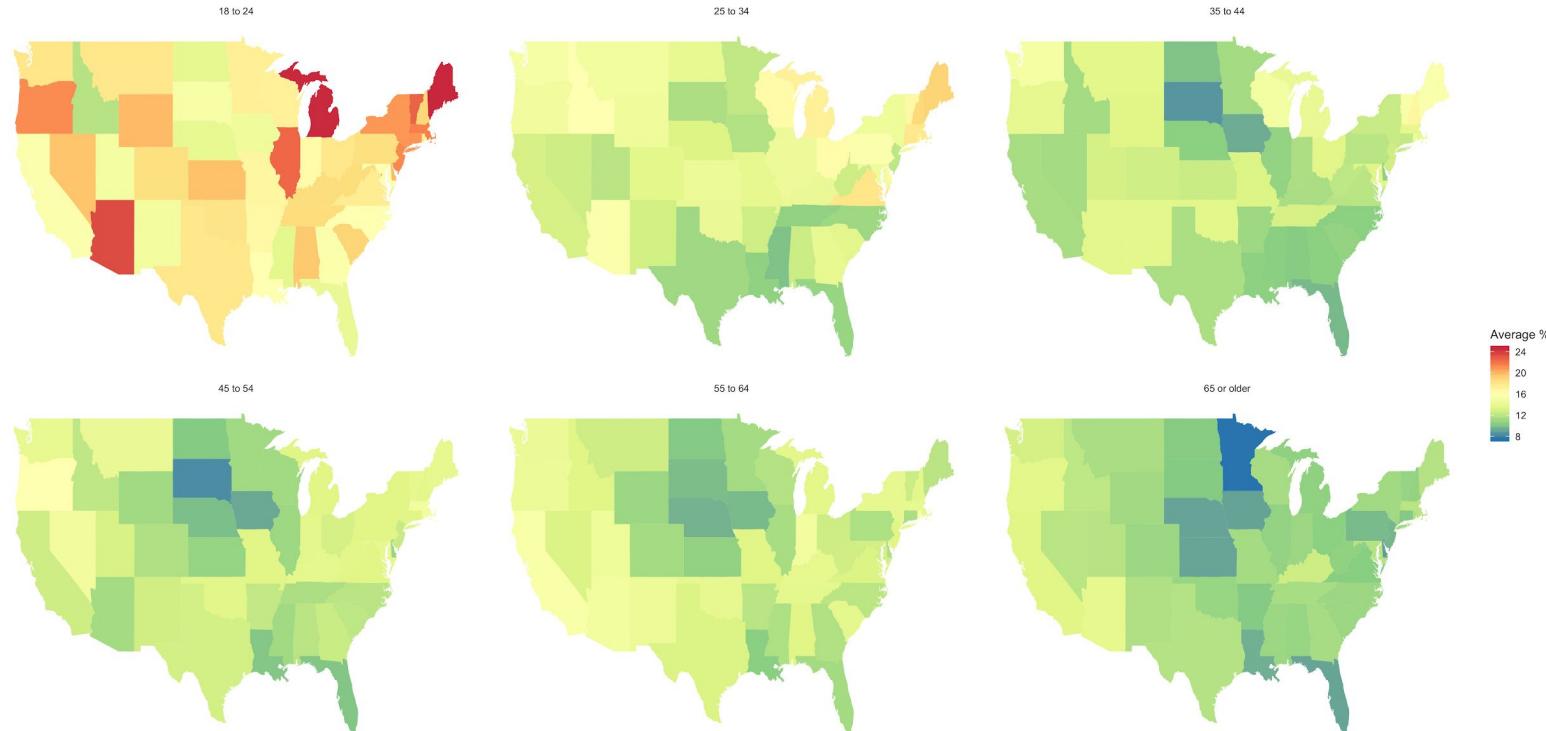


Figure 11. Asthma Prevalence by State and Age Group from 2007 - 2008.

Modeling

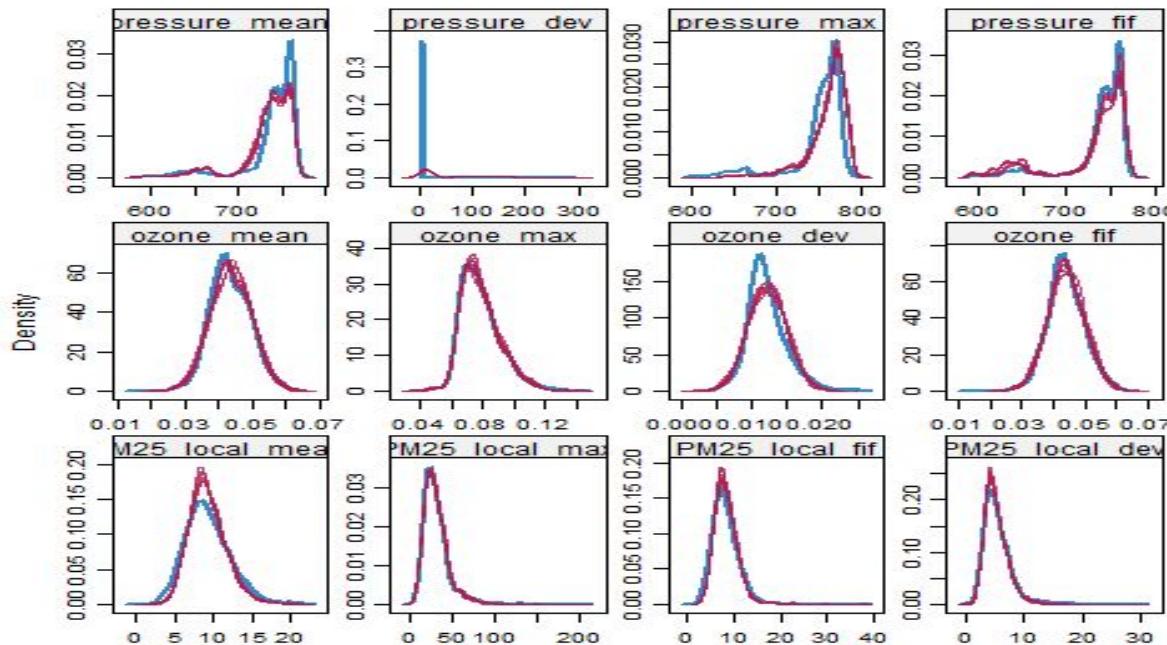
Objective

Model Asthma and COPD health outcomes based on air quality and weather data and identify which of these variables are most important for predicting:

- Asthma Emergency Department Visits
- COPD Hospitalizations
- COPD Mortality

Modeling - Missing Value Imputation

Density Plot of Observed (blue) vs Imputed (red) values.



Missing values were imputed using Conditional Multiple Imputation with the MICE package.

Here blue represents the observed data and red shows the imputed data.

Figure 12. Observed versus imputed values for predictor variables.

Modeling - Variable Importance

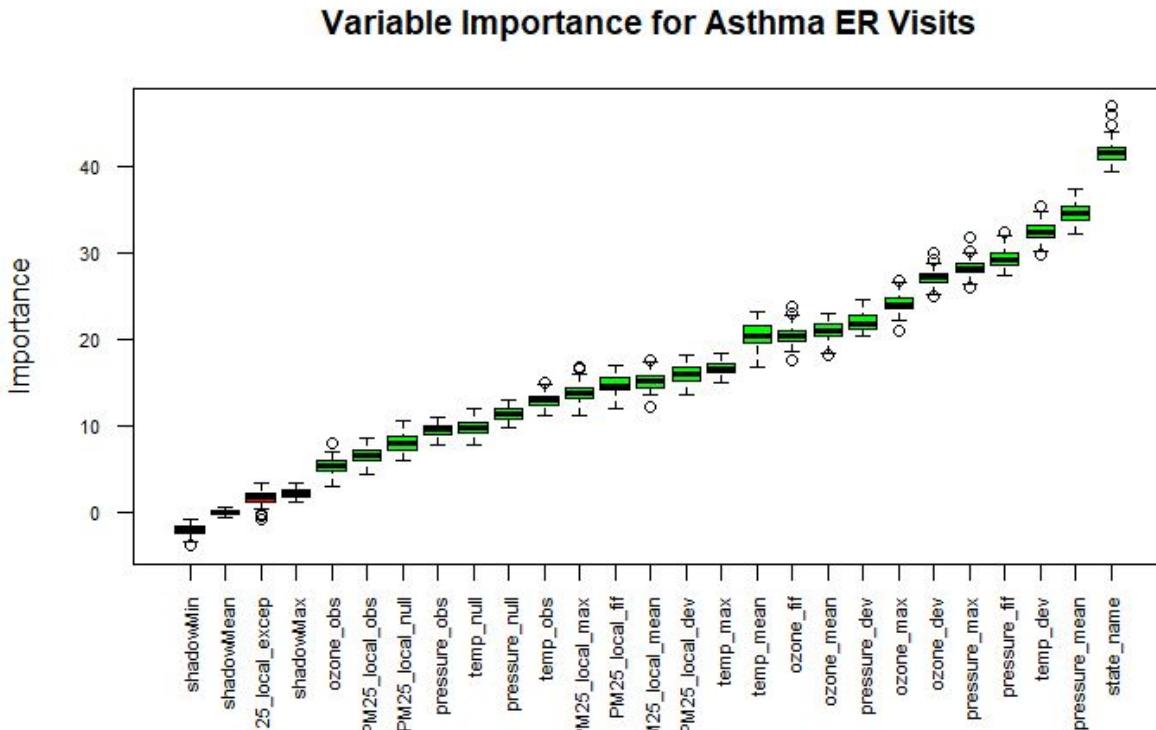


Figure 13. Importance of each predictor variable.

We used Boruta package to rank all variables by importance.

Boruta is a feature ranking and selection algorithm based on a random forests algorithm.

Modeling - Asthma ER Visits

Variable Importance Using Stepwise Backward Selection

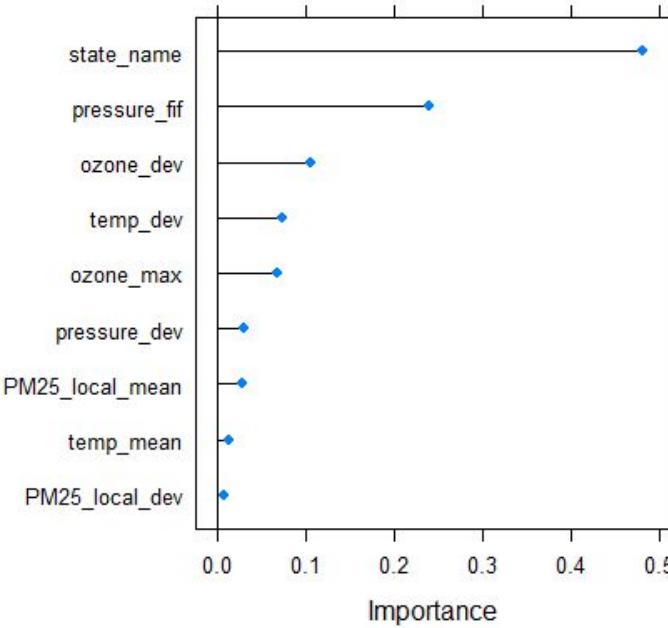


Figure 14. Importance of variables using stepwise backward selection.

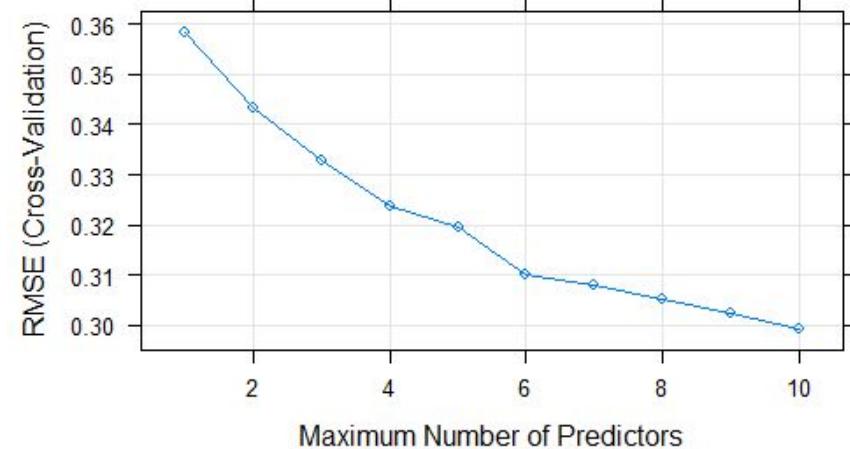


Figure 15. Cross-validated RMSE as number of predictors increases.

RMSE	R-squared	MAE
1.3387169	0.4079990	1.237774

Modeling - COPD Hospitalizations

Variable Importance Using Stepwise Backward Selection

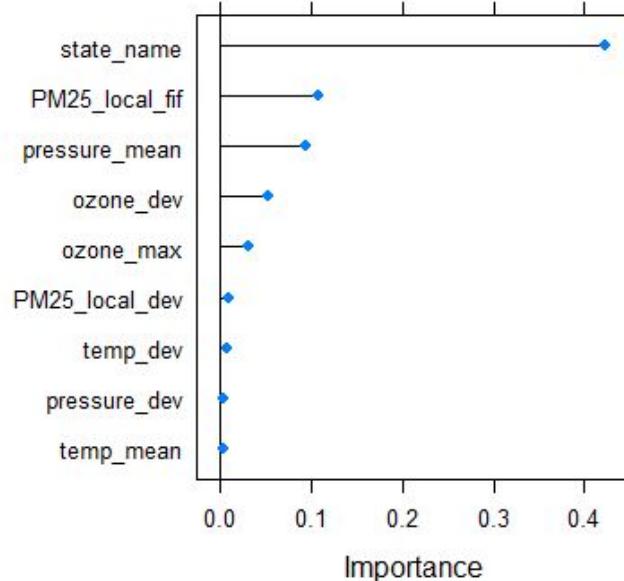


Figure 16. Importance of variables using stepwise backward selection.

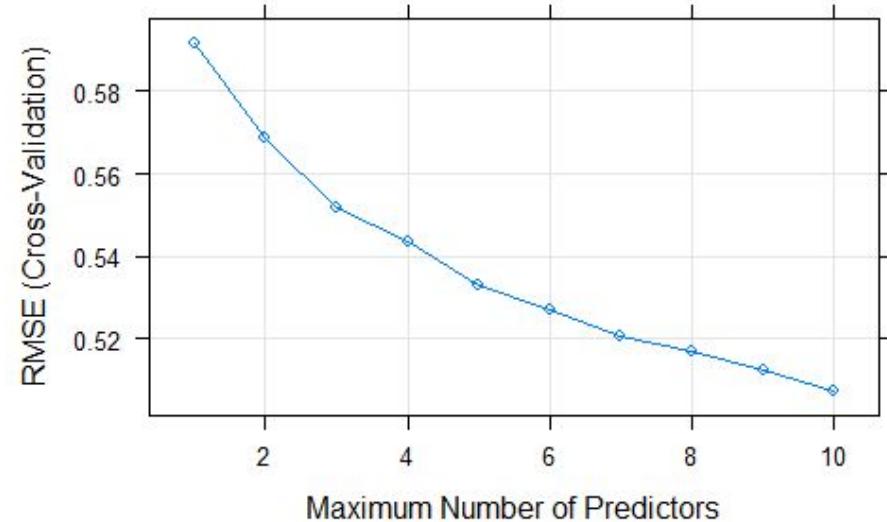


Figure 17. Cross-validated RMSE as number of predictors increases.

RMSE	R-squared	MAE
1.6968100	0.3275164	1.488146

Modeling - COPD Mortality

Variable Importance Using Stepwise Backward Selection

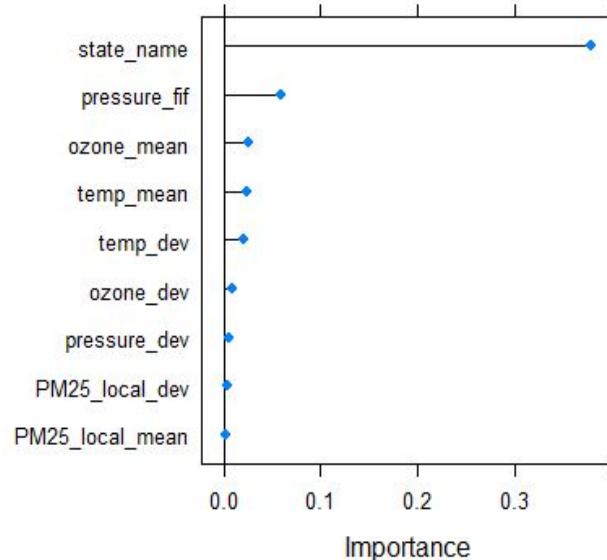


Figure 16. Importance of variables using stepwise backward selection.

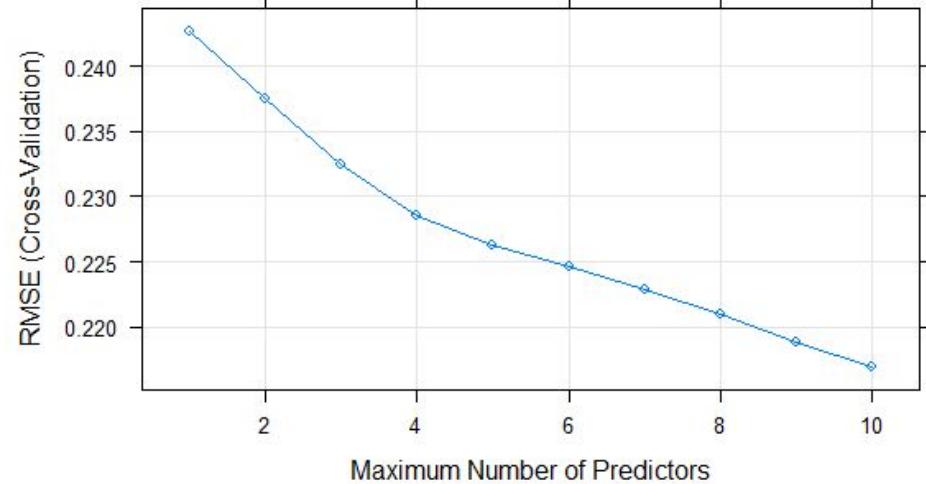


Figure 17. Cross-validated RMSE as number of predictors increases.

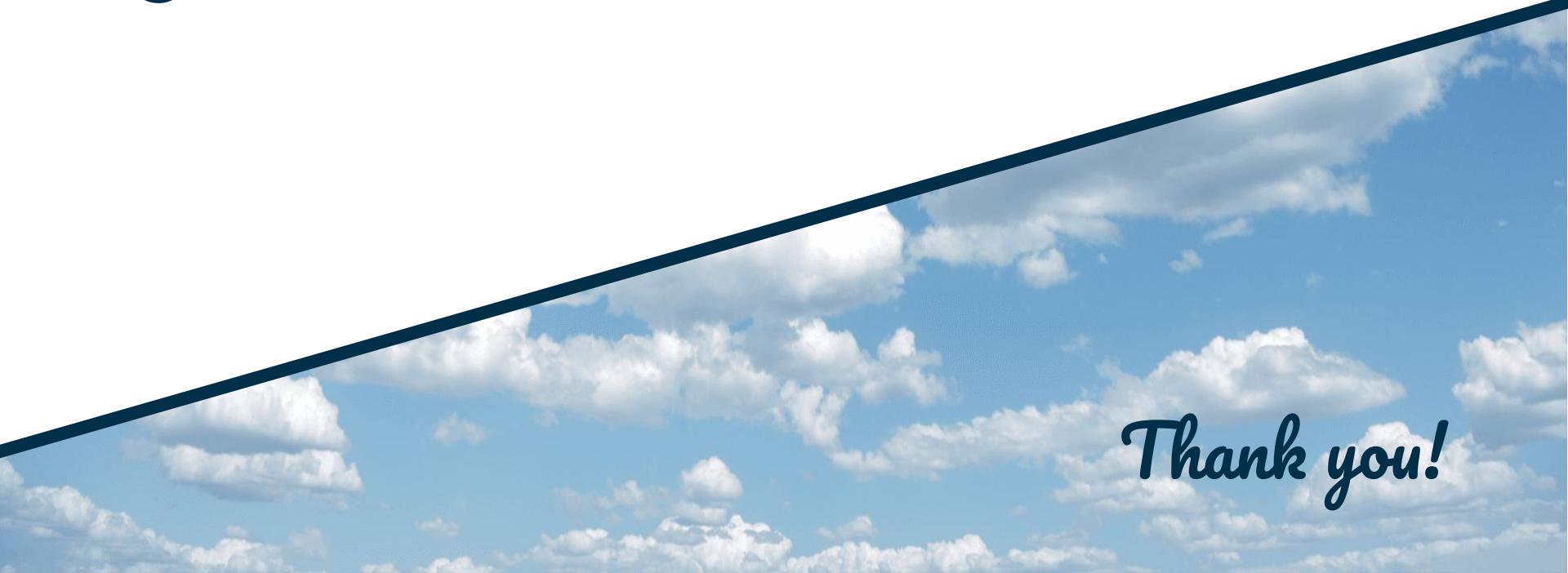
RMSE	R-squared	MAE
1.2450260	0.2021450	1.185401

Results & Conclusions

Top 3 Most Important Variables After Adjusting for State and Normalized RMSE

		RMSE	Rsquared
Asthma ER Visits	1. Ambient Pressure 50th 2. Ozone SD 3. Ambient Temperature SD	1.3387	0.4079
COPD Hospitalizations	1. PM2.5 Local Conditions 50th 2. Mean Ambient Pressure 3. Ozone SD	1.6968	0.3275
COPD Mortality	1. Ambient Pressure 50th 2. Mean Ozone 3. Mean Temperature	1.2450	0.2021

Questions?



Thank you!

References

Environmental Public Health Data Tracking Network, by the Center for Disease Control:

Link to data: <https://ephtracking.cdc.gov/DataExplorer/#/>

Link to background information: <https://www.cdc.gov/nceh/tracking/about.htm>

Historical Air Quality Data from the Air Quality System, by the Environmental Protection Agency:

Link to data: <https://www.kaggle.com/epa/epa-historical-air-quality>

Link to background information: https://aqs.epa.gov/aqsweb/documents/About_aqs_data.html