

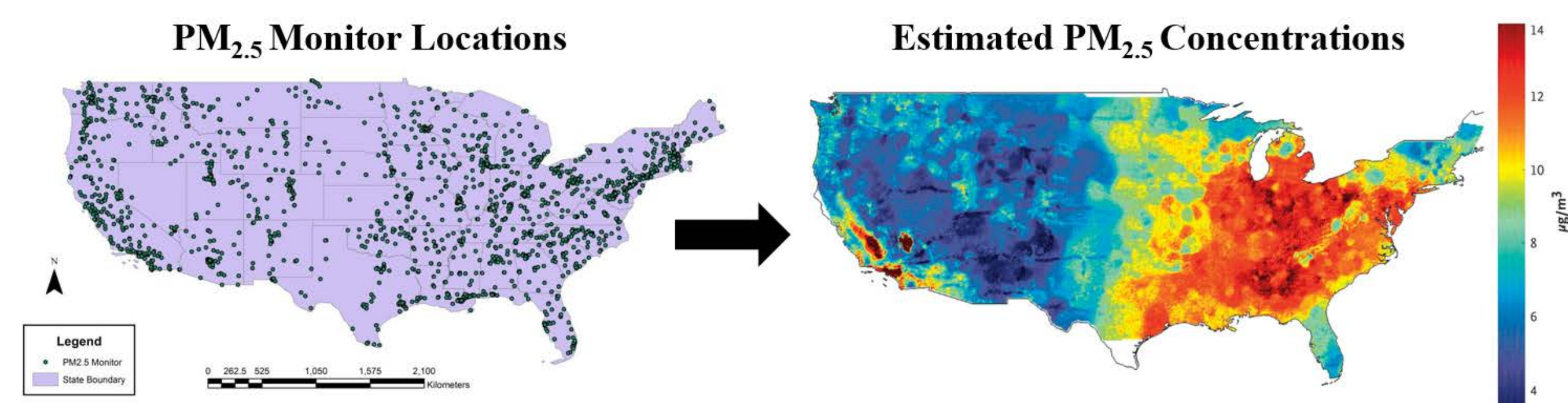
**Members:** Keyan Halperin, Christopher Hase, Justin Lee, Casey Meehan

**Mentor:** David Sondak

**Partner:** Harvard School of Public Health

## Background

PM<sub>2.5</sub> – fine particles with a diameter of 2.5  $\mu\text{m}$  or less – is measured at approximately 2,000 air pollution monitors located throughout the US. However, these monitors are costly to operate and thus, sparsely distributed. Consequently, the air quality is not known for many locations throughout the country. This is problematic for public health studies since we cannot estimate the overall effect that pollution has on health if we do not know what the pollution is in many areas with any degree of certainty.



Thus, using pollution data from these sparsely distributed sensors along with several other types of data, our goal was to fit a model that accurately interpolates pollution throughout the US. With this, it would be possible to create of a continuous map of pollution on the daily level, which could be used by researchers to establish causal relationships between pollution and health outcomes, among several other potential research applications.

## Data

### Monitor Data

Pollution readings from approximately 2,000 PM<sub>2.5</sub> sensors each day for the past 16 years.

### Satellite Data

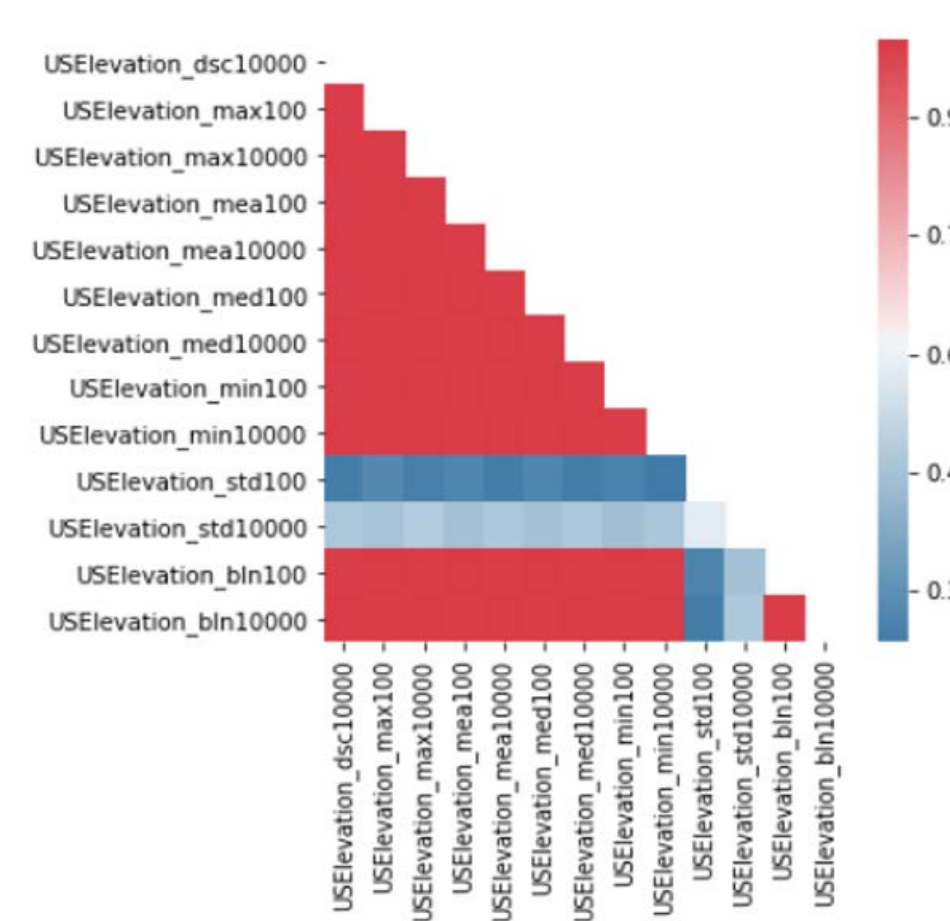
Various meteorological, geographical, and topographical variables collected from satellites and other instruments.

### US Census Data

Demographic variables on the ZIP Code level as collected by the US Census.

Because of the high number of variables, we decided to perform variable selection. Having extraneous or redundant variables not only increases the computation time, but it can also harm the predictive performance of machine learning models. To reduce dimensionality, for each set of highly correlated features, we removed all but one. We also removed variables that had non-significant partial correlations with PM<sub>2.5</sub> after controlling for nearby PM<sub>2.5</sub>. On the right you can see a set of variables for which there were many high pairwise correlations.

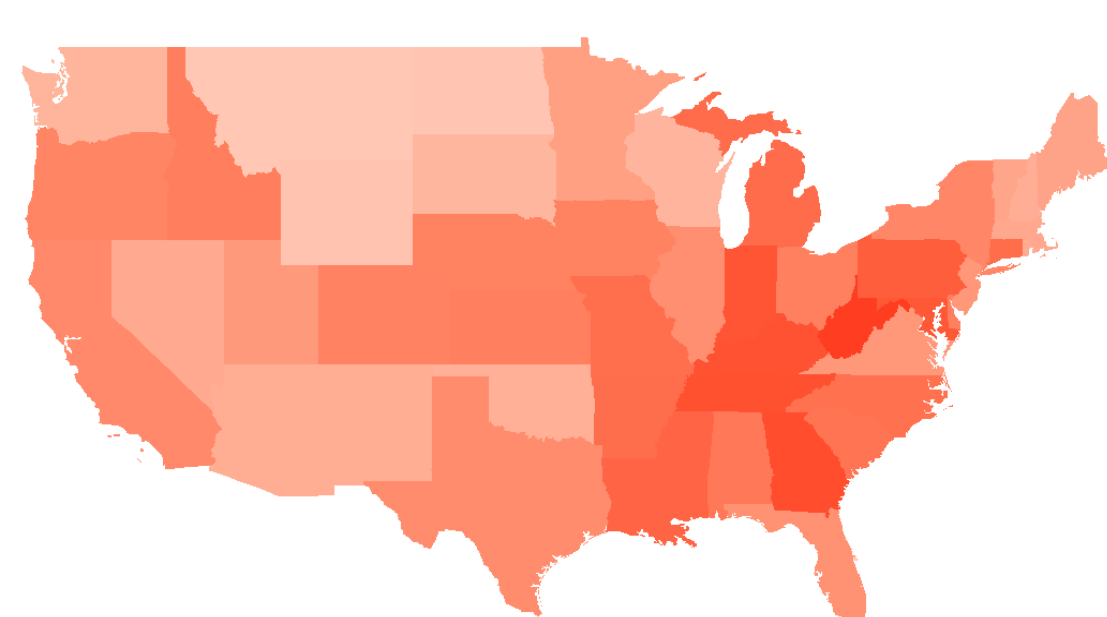
Elevation Pairwise Correlations



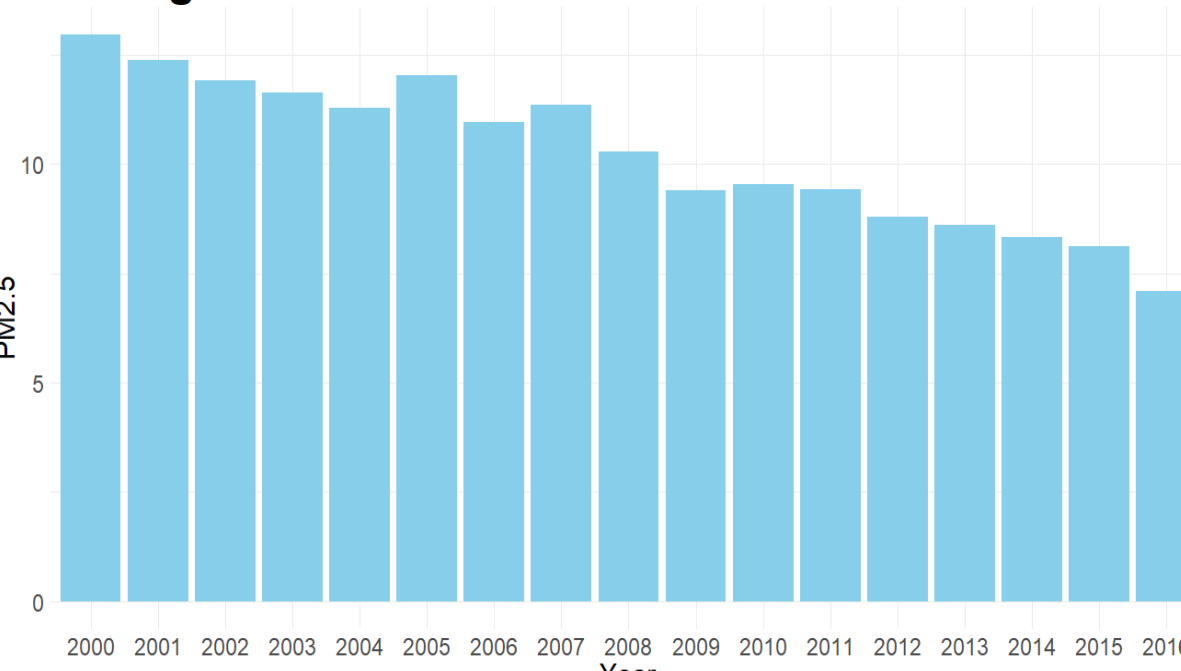
## Exploratory Data Analysis

We also conducted a thorough exploratory data analysis to gain insights that would help us in modeling. Most notably, we discovered that PM<sub>2.5</sub> has strong spatial and temporal dependencies. Below are plots of average PM<sub>2.5</sub> levels by state and average PM<sub>2.5</sub> levels over time that demonstrate this.

Average Pollution by State



Average Pollution over Time

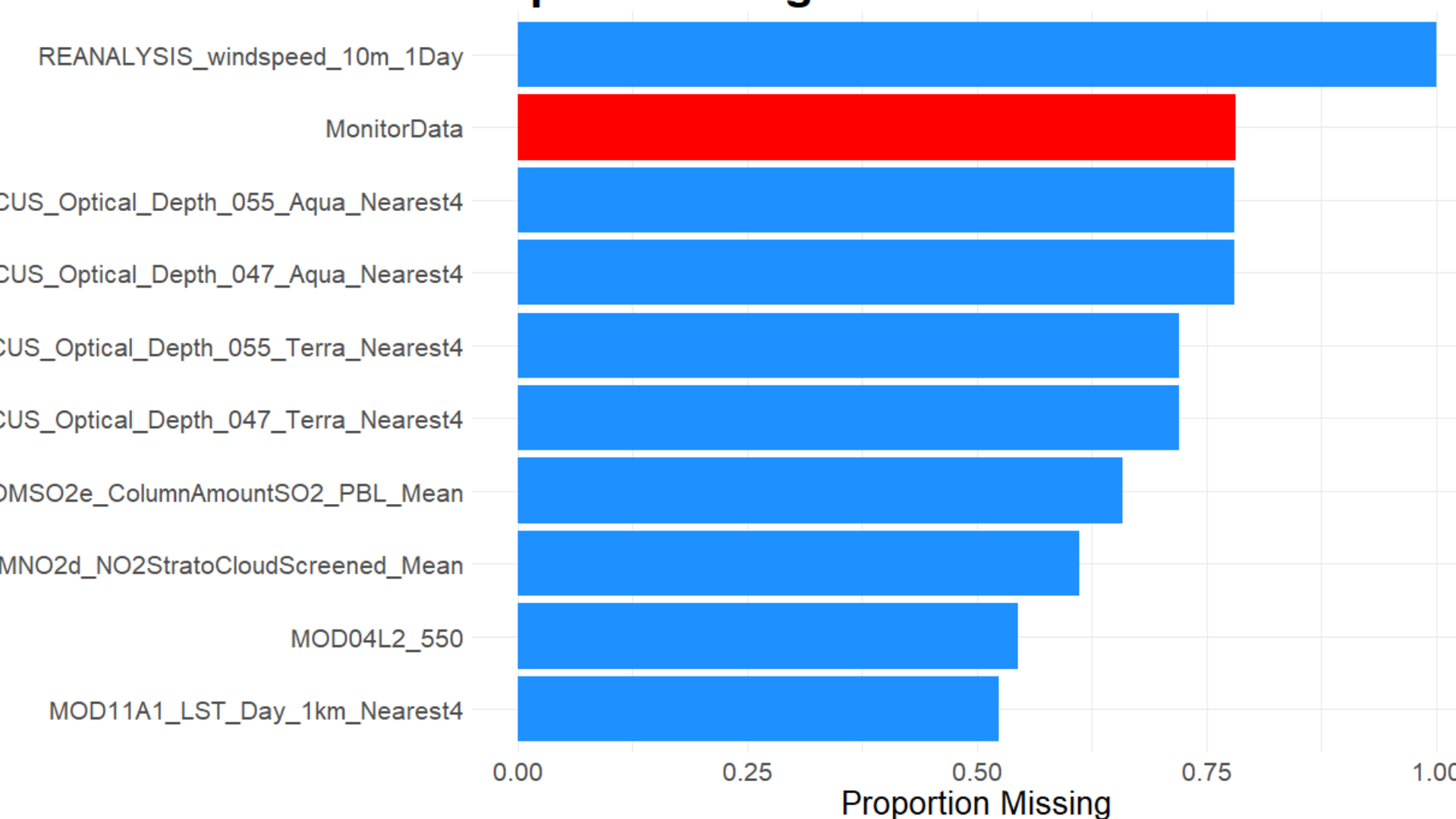


We also explored the relationship between each feature and PM<sub>2.5</sub>. By far, the feature that had the highest correlation with PM<sub>2.5</sub> was nearby PM<sub>2.5</sub>, with a correlation of about 0.857. Thus, there is strong evidence to suggest that we should implement models that take into account both spatial and temporal relationships.

## Imputing Missing Data

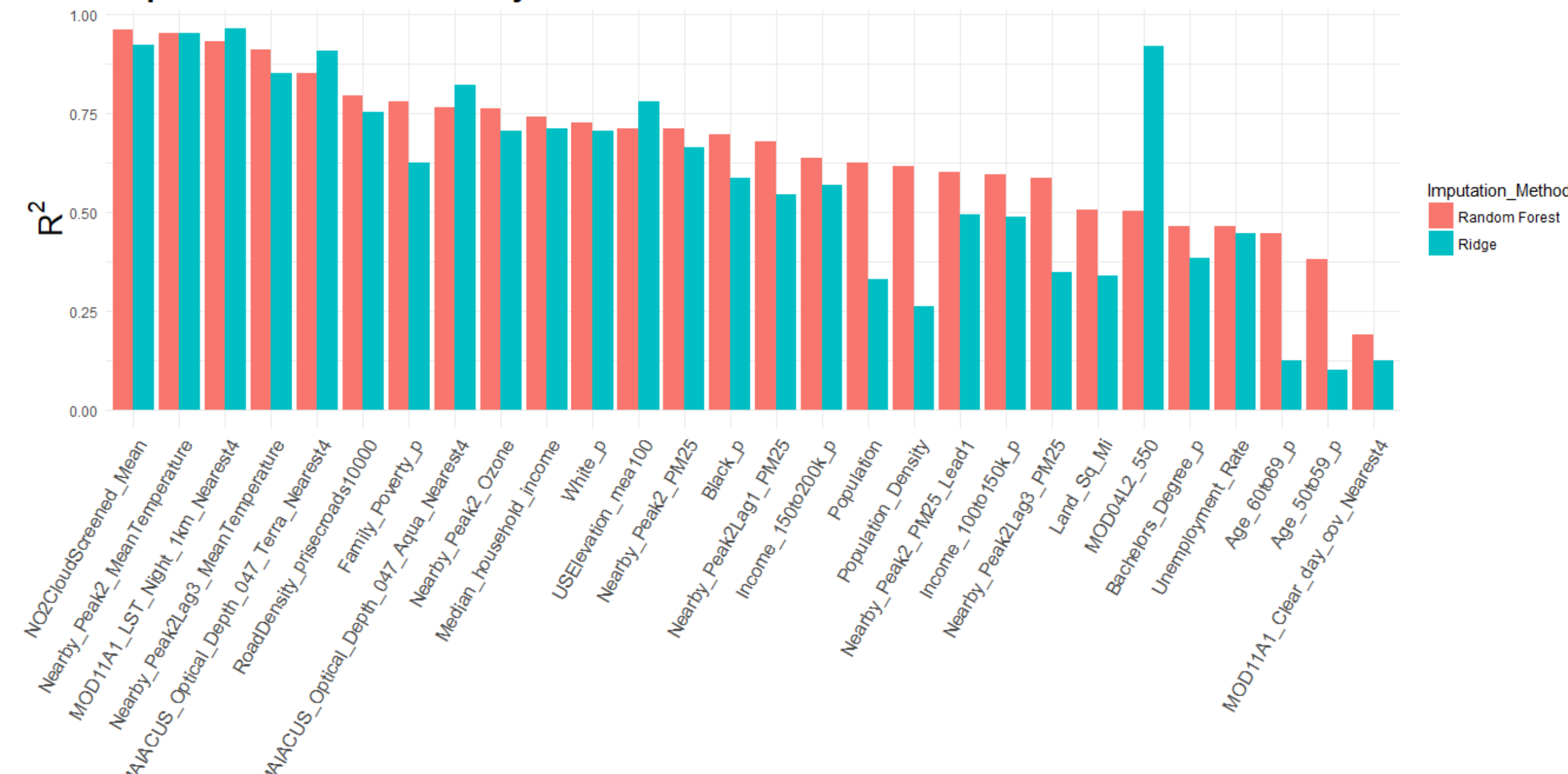
A major obstacle that we had to overcome was the amount of missing data. Many of the variables and almost all observations had some missing data. Even PM<sub>2.5</sub> itself was often missing since the monitors do not measure pollution each day. Below is a graph of the top 10 most missing variables with PM<sub>2.5</sub> monitor data highlighted in red.

Top 10 Missing Variables



Because of the volume of data that was missing and its potential importance for modeling PM<sub>2.5</sub>, we could not afford to ignore it. After some research and experimentation, we decided to implement an iterative imputation algorithm called MissForest (Stekhoven & Bühlmann, 2012). Allowing for other supervised learning models to be used instead of random forest enabled us to compare two MissForest variations – one that used random forests as in the original paper and one that used ridge regressions. The scheme we created for evaluating the quality of our imputations indicated that our imputations were accurate. The random forest variation achieved a weighted R<sup>2</sup> of 0.717 and the ridge variation achieved a weighted R<sup>2</sup> of 0.774. Below is a plot of the imputation performance of both variations of this algorithm for each variable as measured by R<sup>2</sup>.

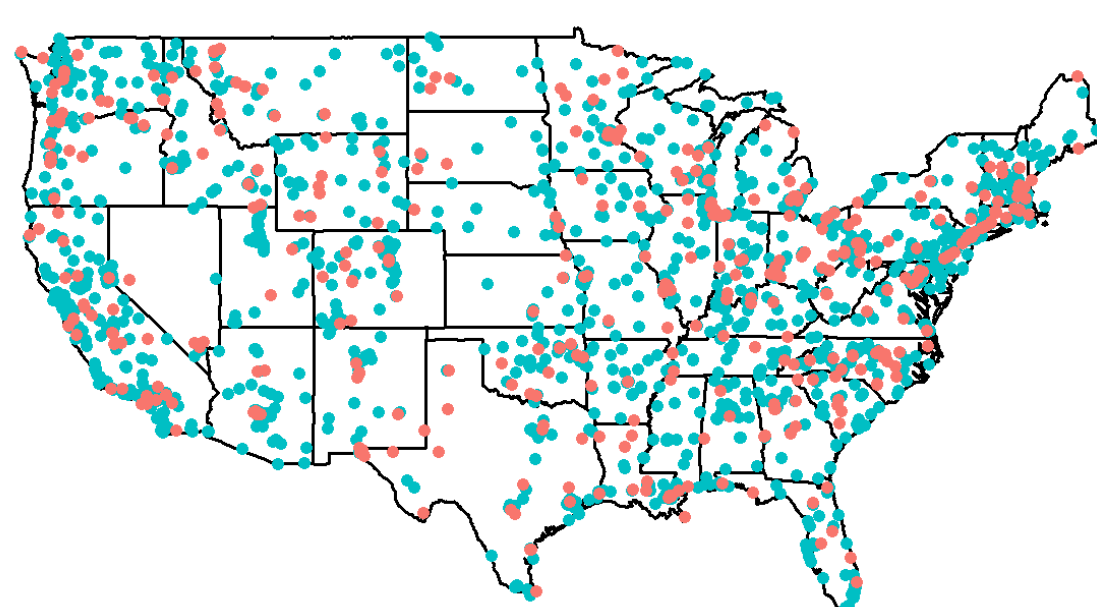
Imputation Performance by Variable



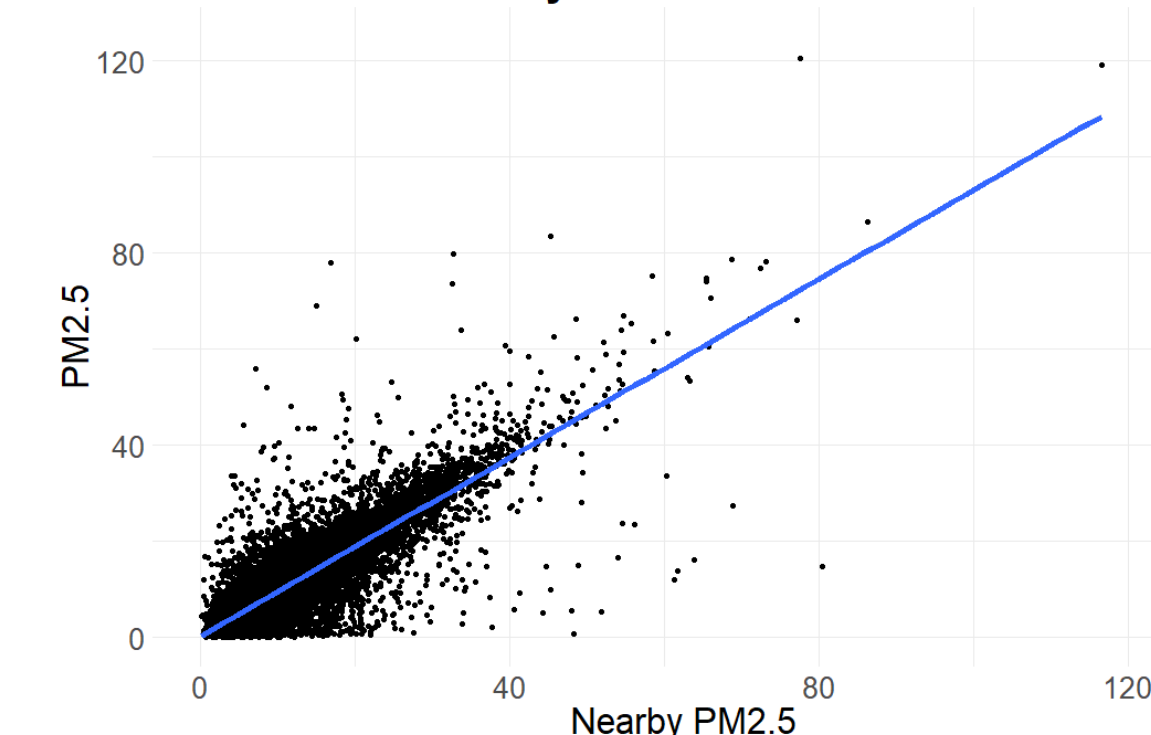
## Baseline Model

After pre-processing and imputation, we tuned our models using K-fold cross-validation on a random 80% of the sensor site sequences and reserved the remaining sensors to test the performance of our models.

Sensor Train-Test Split

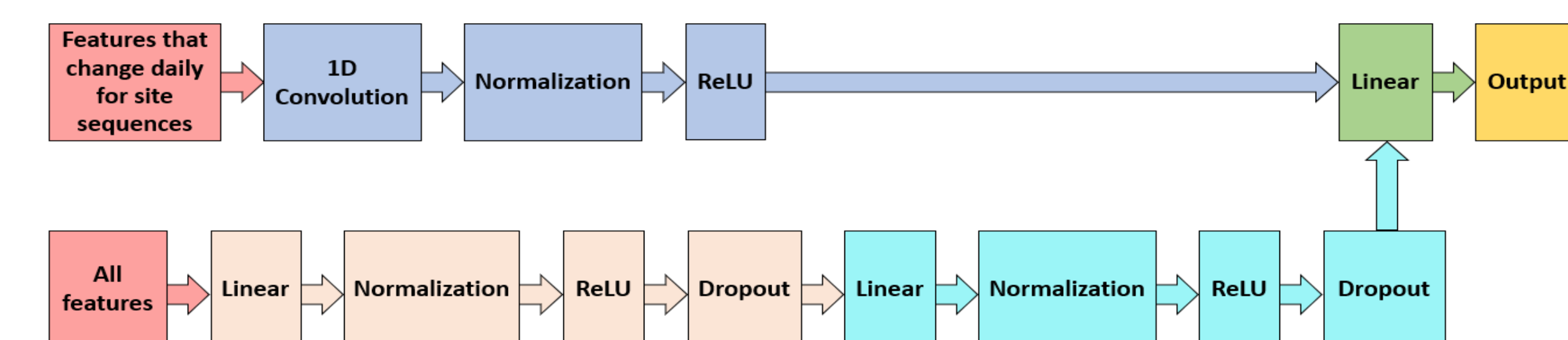


Pollution vs. Nearby Pollution



Because of the very high correlation between PM<sub>2.5</sub> and nearby PM<sub>2.5</sub>, we decided that our baseline model should be a simple linear regression with nearby PM<sub>2.5</sub> as the only feature. Despite its simplicity, this model performs surprisingly well with a test R<sup>2</sup> of 0.712. Above is a scatterplot of the relationship between PM<sub>2.5</sub> and nearby PM<sub>2.5</sub> on a random 1% subset of the data.

## Modeling

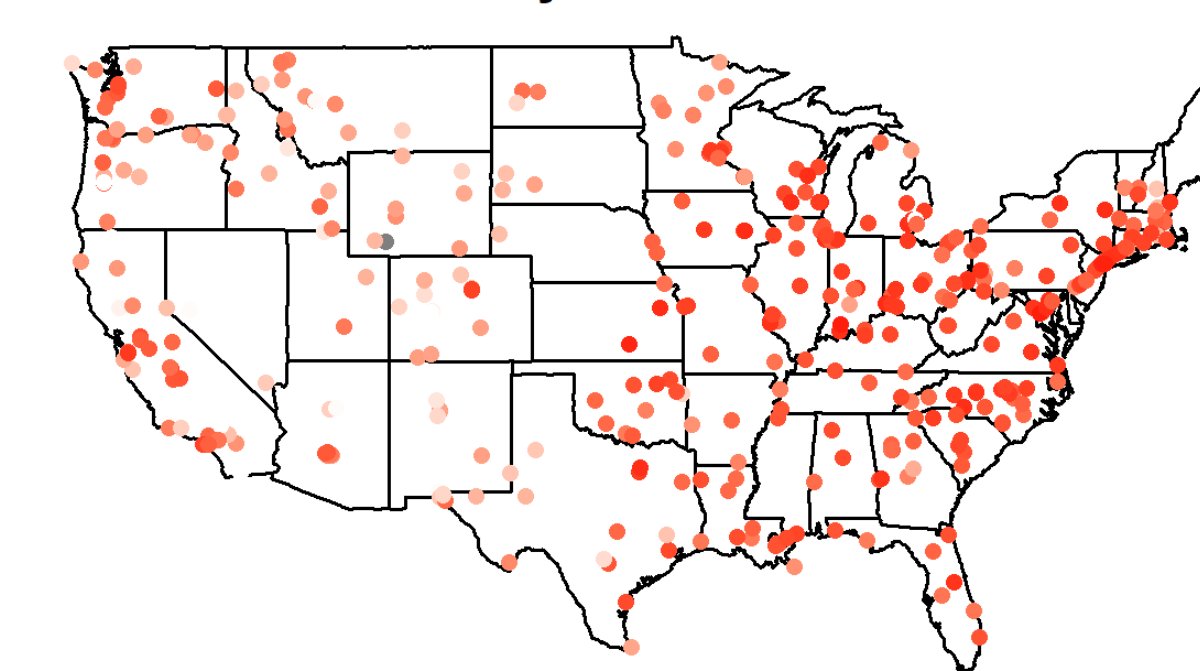


Because consecutive days are likely to be related in ways that are relevant for predicting pollution, we decided to implement a CNN. In particular, we used a kernel width of size 3 so that features from the previous day, current day, and the following day are used for predicting pollution for the current day. We are still in the process of tuning the CNN, but our best result so far is an R<sup>2</sup> of 0.775.

	OLS	Ridge	RF	XGBoost	CNN	Ensemble
R <sup>2</sup>	0.712	0.733	0.780	0.776	0.775	0.784

We also implemented various standard machine learning models to compare their performances. Although random forest outperformed the other models, our best results came from an ensemble of random forest, XGBoost, and CNN.

Model Performance by Location



On the left is the performance of our ensemble by location as measured by R<sup>2</sup>. Unsurprisingly, the ensemble produces much more accurate predictions in regions with more sensors. This highlights the need to install more monitors in locations where monitors are sparse.

## Discussion

Because of the disproportionate importance of nearby PM<sub>2.5</sub>, we believe that it is absolutely essential for more pollution monitors to be installed, especially in regions where there are few. We also have evidence to suggest that our imputation procedure provides high quality imputations, and since the procedure is quite easy to implement, we recommend that HSPH use it in the future.

Going forward, in addition to continuing to improve model performance, we believe that an important task is to quantify the uncertainty of the PM<sub>2.5</sub> predictions. This would be important because (1) it would give more insight into model performance in areas that are far away from sensors (2) it would allow for more accurate variance estimates of any associated causal effects and (3) it could be used in determining which locations should be prioritized for new sensor installations. Gaussian process regression may be a good tool for quantifying the prediction uncertainty.

## Software Package

The scripts written for this project were all run on Harvard's Odyssey Research Computing Cluster, and instructions for running our code were written with Odyssey in mind. This is because of the substantial amount of RAM and processing power needed to load and fit models on the nearly 20 GB of pollution data. Python and R scripts for pre-processing and imputation as well as model validation, training, and testing were wrapped in Bash scripts formatted to be scheduled using Slurm. All code was made available on a GitHub repository so that we could turnover our work to HSPH for use and continued development.

