

# What Can We Learn From Machine Learning Models Developed for Short-term Forecasting of PM<sub>2.5</sub>?

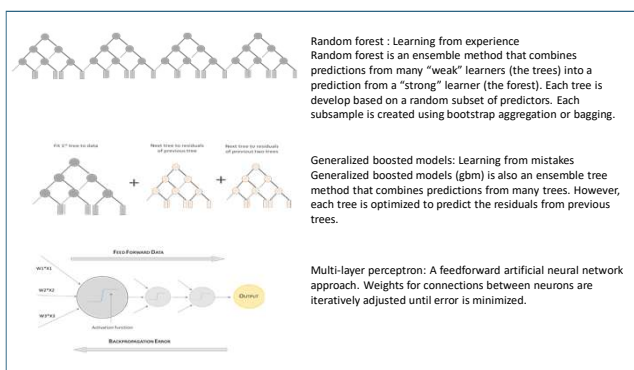
Meenakshi Rao<sup>1,2,\*</sup>, Aaron Fellows<sup>2</sup>, Morgan Schafer<sup>2</sup>, Phillip Orlando<sup>1</sup>, and Linda Acha George<sup>1</sup>

<sup>1</sup>Portland State University, Portland OR, United States, <sup>2</sup>Oregon Department of Environmental Quality, Portland, OR, United States

\* Corresponding author mrao@pdx.edu

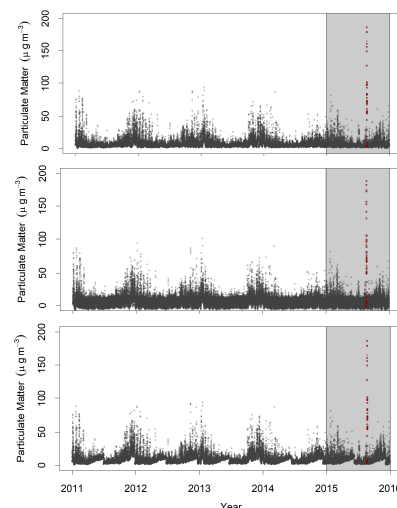
## Background

Accurate short-term forecasting of PM<sub>2.5</sub> concentrations is important for reducing exposure to unhealthy levels of air pollution. The EPA uses a 12-hour weighted mean for its NowCast prediction of PM<sub>2.5</sub> concentrations. We compared the performance of three Machine Learning (ML) algorithms (random forest, generalized boosted trees and multi-layer perceptron) against the EPA's NowCast.



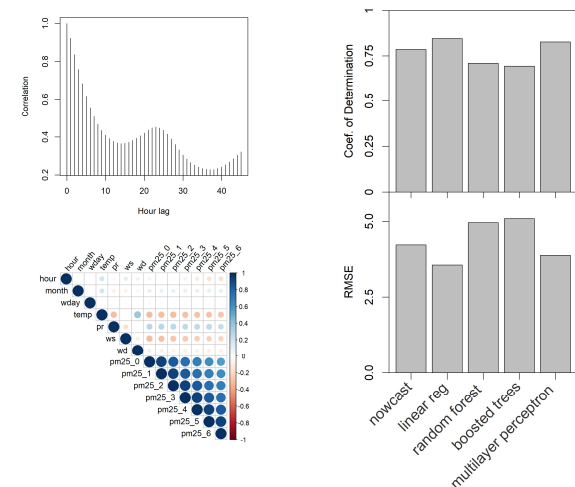
Low-cost air pollution sensors are becoming more prevalent, and provide additional opportunities for PM<sub>2.5</sub> forecasting. However data from these sensors may have more noise and drift compared to traditional higher-end instruments. Here we examine how the EPA's NowCast and ML algorithms perform on data that have been perturbed to include noise and drift. The NowCast method was especially developed to capture rapid changes in PM<sub>2.5</sub> associated with a wildfire or smoke event. Therefore we also examine how these algorithms perform in predicting PM<sub>2.5</sub> during a wildfire event.

## Data Input for Machine Learning (ML) Algorithms



PM<sub>2.5</sub> from Oregon DEQ's Portland air monitoring station was used for model inputs. PM<sub>2.5</sub> from 2011-2014 was used to train the algorithms. PM<sub>2.5</sub> from 2015 (grey background) was used for validation. Wildfire smoke created a spike in PM<sub>2.5</sub> during 2015, and is shown in red. The top panel shows unperturbed data. The middle panel is unperturbed data with  $\pm 5$  standard deviation of added noise. The bottom panel shows unperturbed data with added drift.

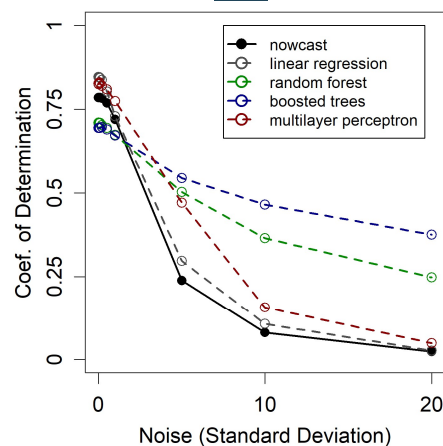
## How do ML Algorithms Perform for Short Term Air Forecasting?



PM<sub>2.5</sub> was first examined without noise or drift (above). The upper left panel shows autocorrelation in unperturbed PM<sub>2.5</sub>, indicating hourly and diurnal patterns. The lower left panel shows correlation for all predictors used in the machine learning models. The right panel shows the performance of the EPA's NowCast, linear regression, random forest, boosted trees, and multilayer perceptron approaches. Coefficient of determination ( $R^2$ ) and Root Mean Square Error (RMSE) were used to assess performance, and were determined by comparing predicted and observed PM<sub>2.5</sub> during the 2015 validation period.

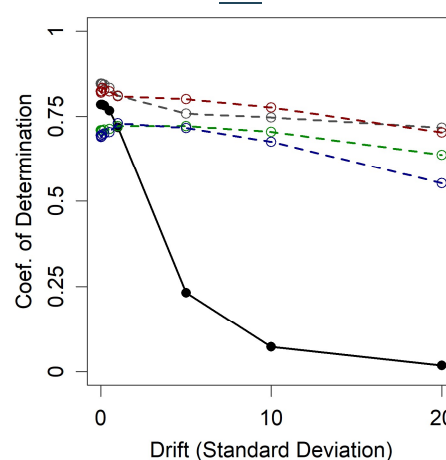
## How do ML Algorithms Perform with Noise, Drift, and Outliers?

### Noise



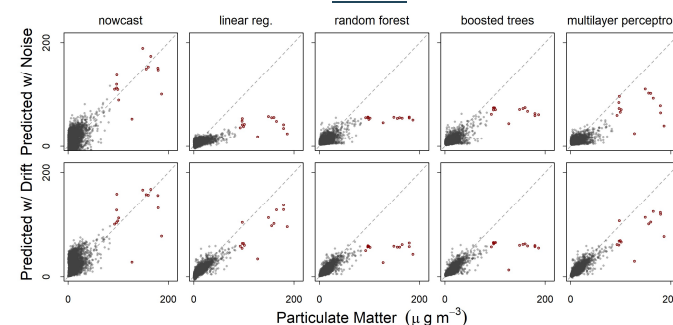
The coefficient of determination ( $R^2$ ) between predicted and observed 2015 validation data for model runs with added noise. Model performance remained high at small-to-moderate noise.  $R^2$  dropped with large increases in noise for all models. Random forest and boosted trees performed best at high noise.

### Drift



The coefficient of determination ( $R^2$ ) between predicted and observed 2015 validation data for model runs with added drift. Performance of all algorithms (with the exception of the NowCast) remained relatively stable with increased drift.

### Outliers



Predicted vs. Observed PM<sub>2.5</sub> ( $\mu\text{g m}^{-3}$ ) with  $\pm 5$  S.D. of added noise and drift for all model runs. Smoke from a fire created a spike in PM<sub>2.5</sub> during Aug. 2015. Capturing these PM<sub>2.5</sub> spikes is particularly important for short-term air quality forecasting. The EPA NowCast and multilayer perceptron models performed well, whereas other models under predicted PM<sub>2.5</sub>. These figures also show several models often under-predicted PM<sub>2.5</sub> despite high  $R^2$ .

## Conclusions

ML algorithms performed well in the presence of noise and drift, although they did not perform as well on predicting values outside the range of the training datasets (outliers). Overall, ML algorithms show promise in generating good predictions based on relatively noisy datasets characteristic of low-cost sensor networks.