



The Application of Bagged Decision Tree to Corrections of Wave Height Forecasts

Ashley Ellenson, Yuanli Pei, Greg Wilson, Tuba Ozkan-Haller, Xiaoli Fern
OSU Coastal and Ocean Engineering Brown Bag

April 19, 2018



◀ Hide

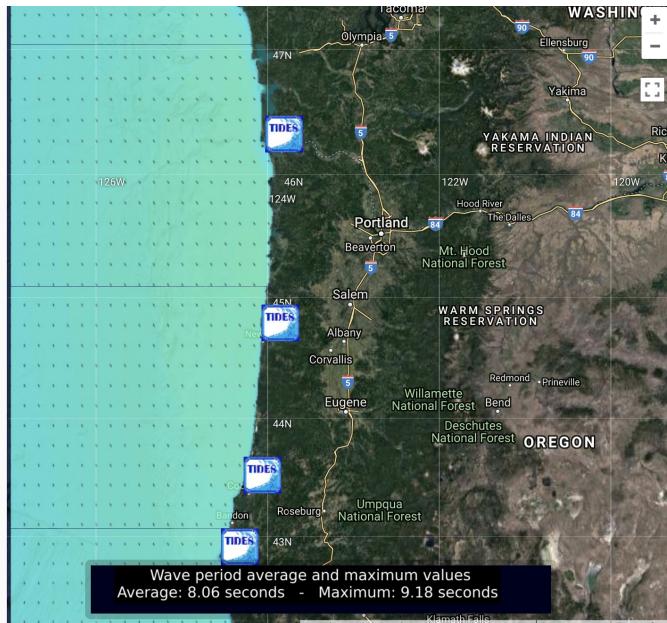
The keys for Currents and Wave Direction are estimates. See About page for details.

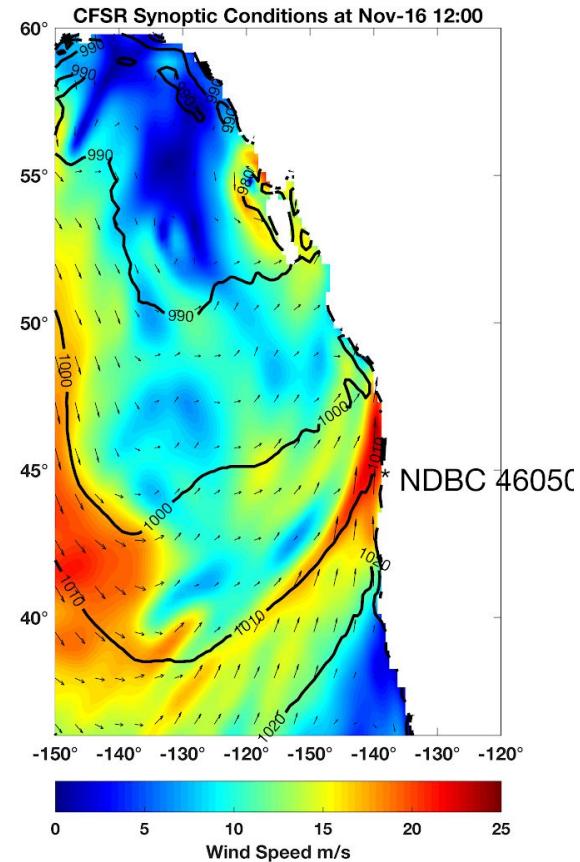
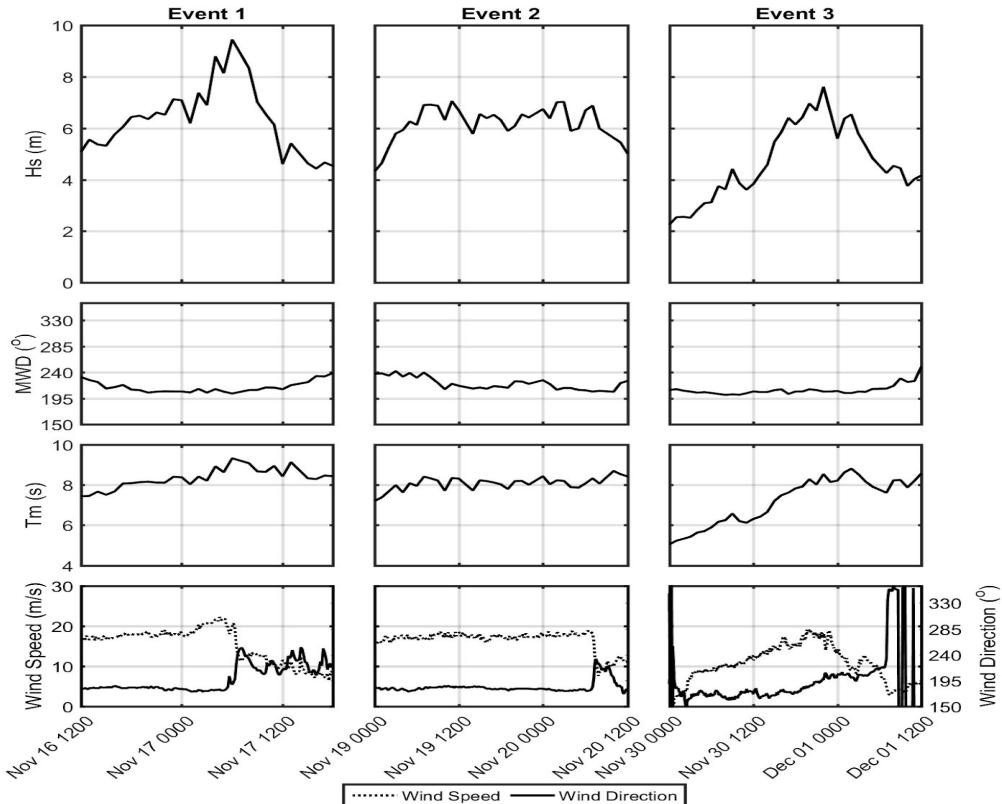
Forecast Date/Time

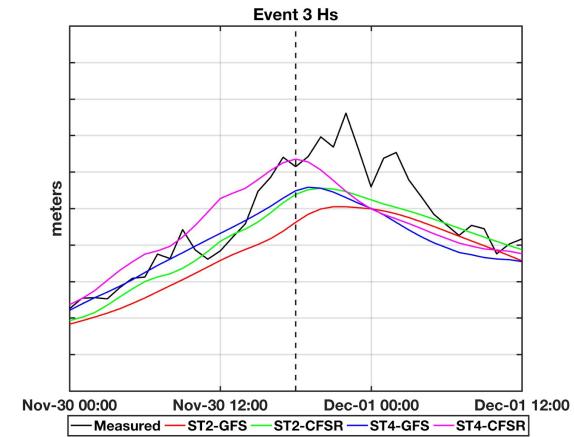
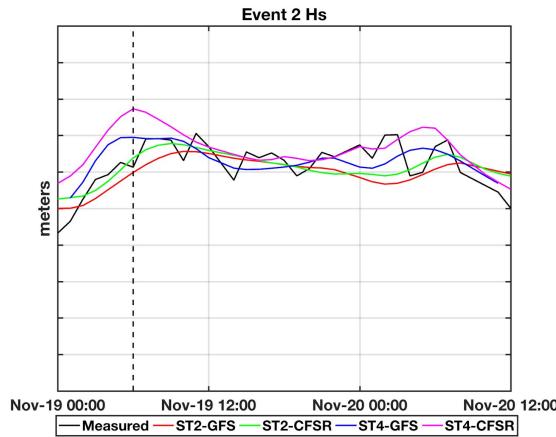
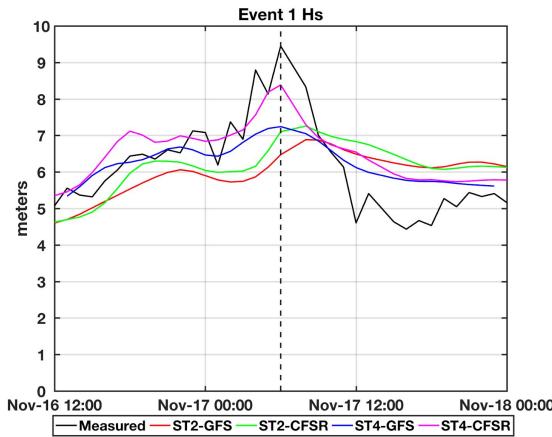
Prev April 17, 2018, 5 p.m. Next

Map Data Layers

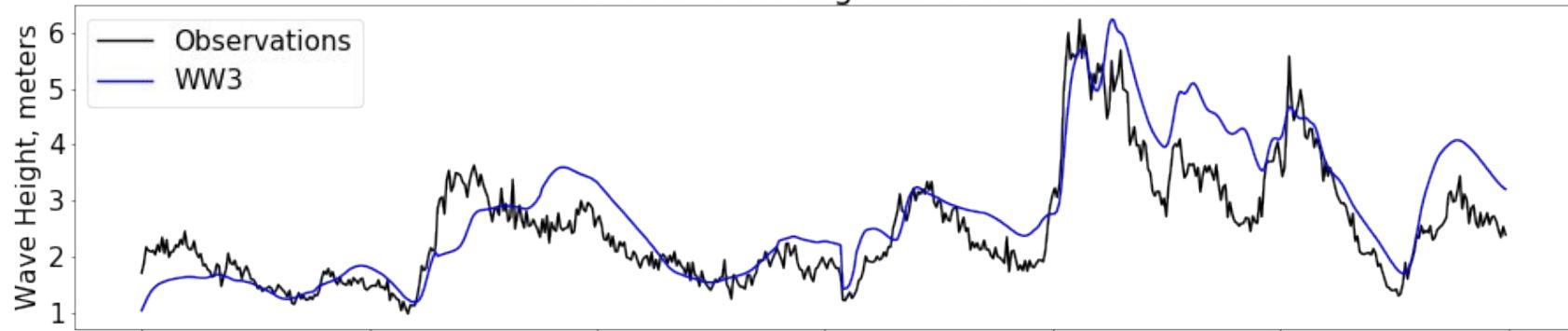
- Surface Temperature OFF
- Surface Currents OFF
- Wave Height ON
- Wave Direction/Period ON
- Wind OFF
- Bottom Temperature OFF
- Surface Salinity OFF
- Bottom Salinity OFF
- Sea Surface Height OFF
- Thermocline OFF



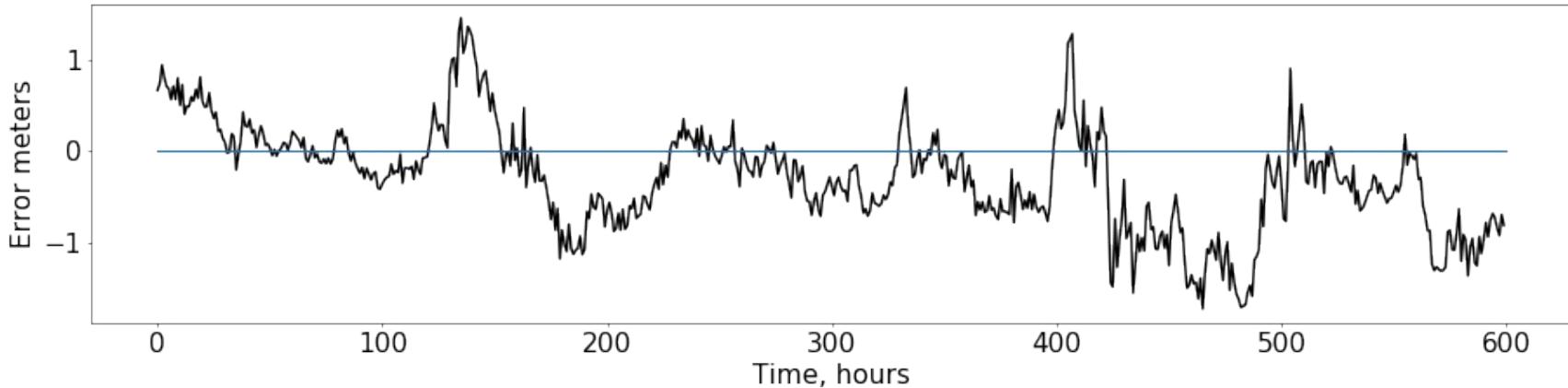




Winter Wave Height Time Series



Winter Error Time Series

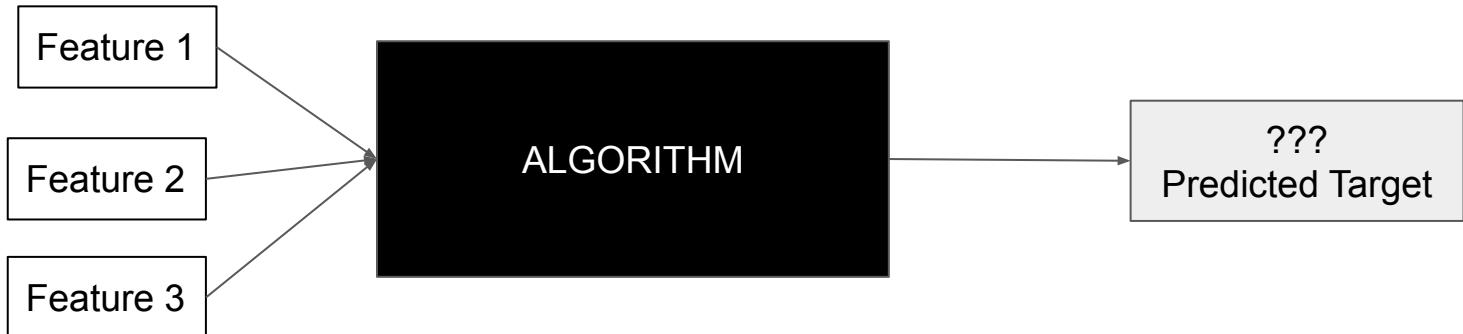


Machine Learning Algorithms: A General Overview

TRAIN



TEST



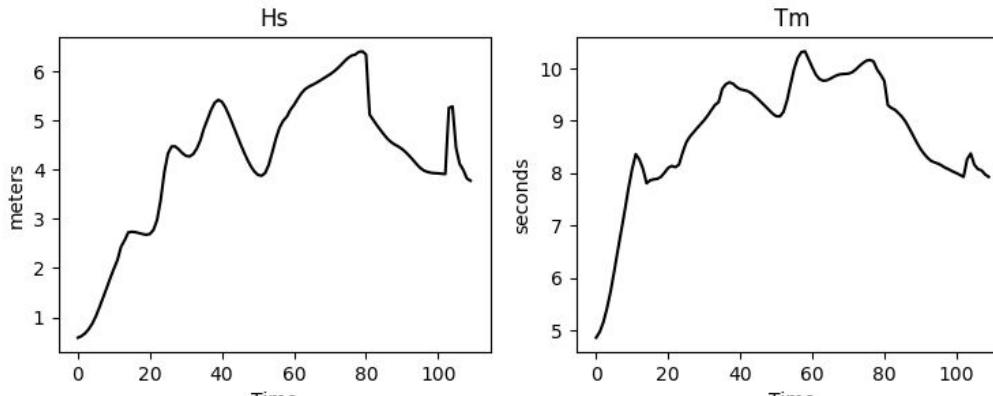
Decision Tree On Wave Forecasts: A Simple Case



Decision Tree On Wave Forecasts

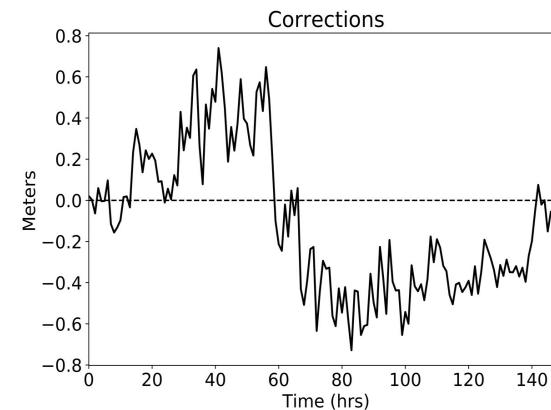
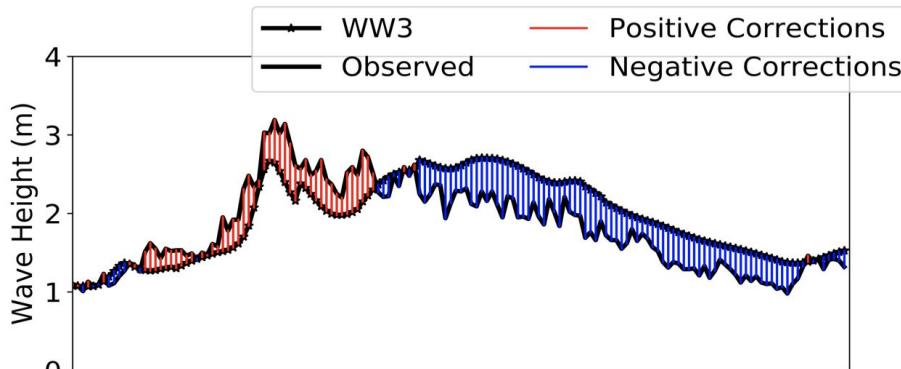
Input Features

$$x(t) = [H_s(t), T_m(t)]$$

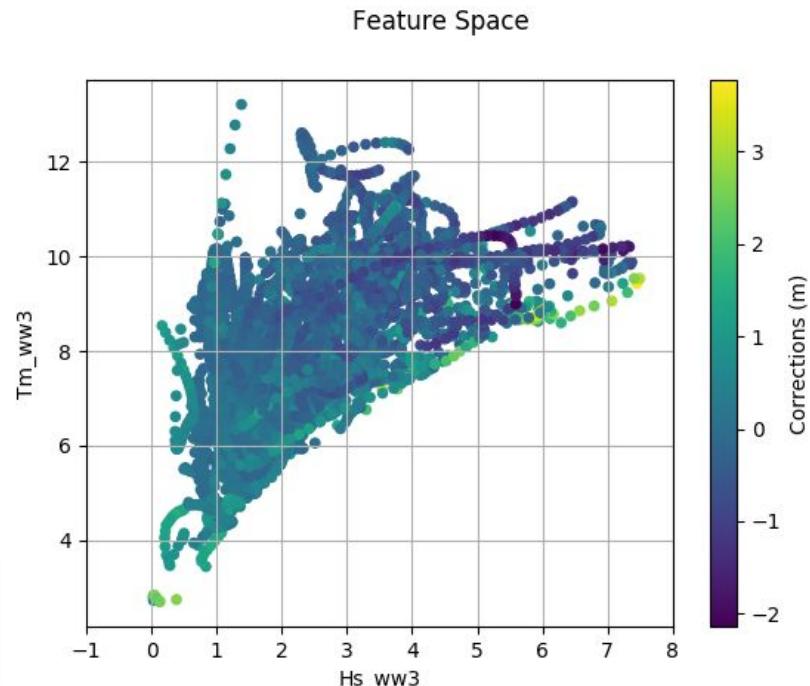
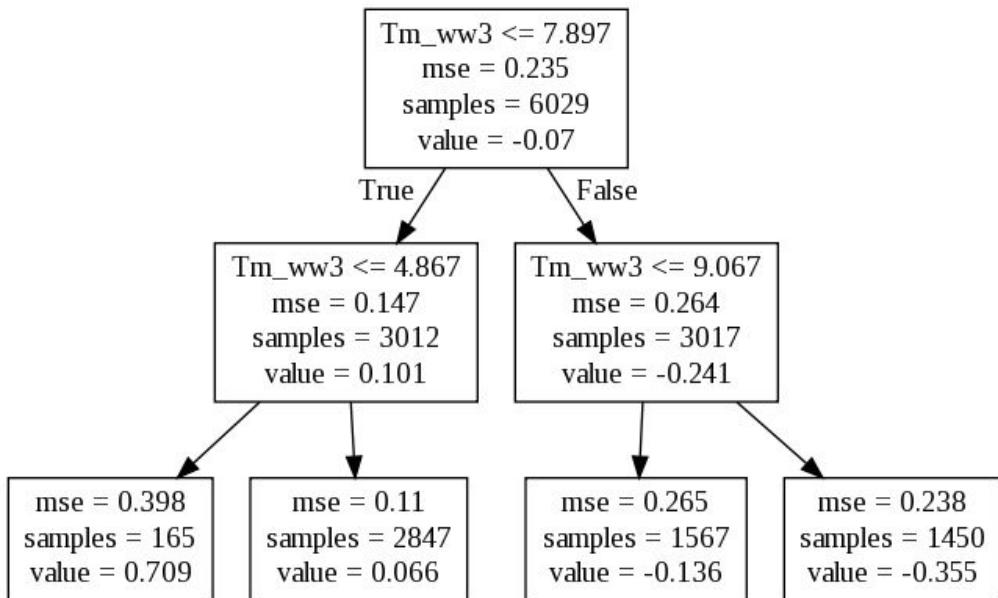


Target

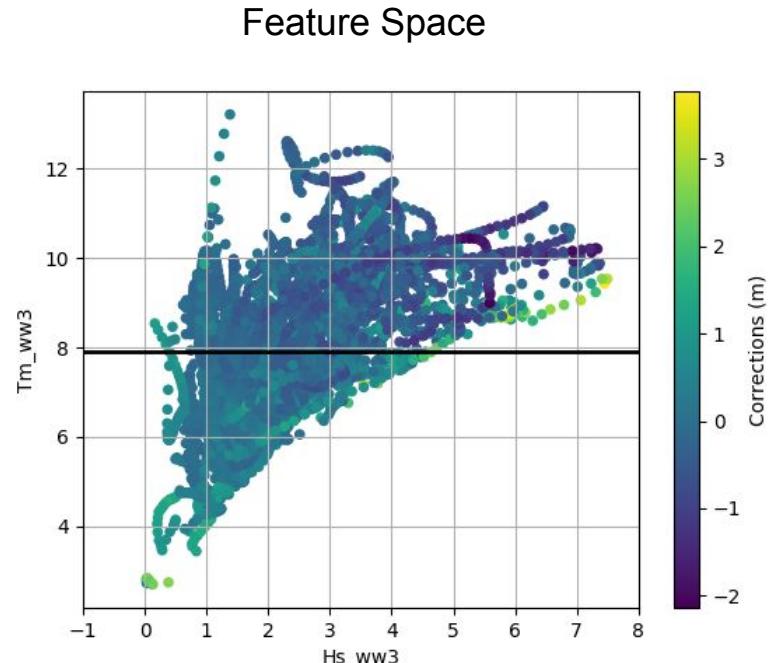
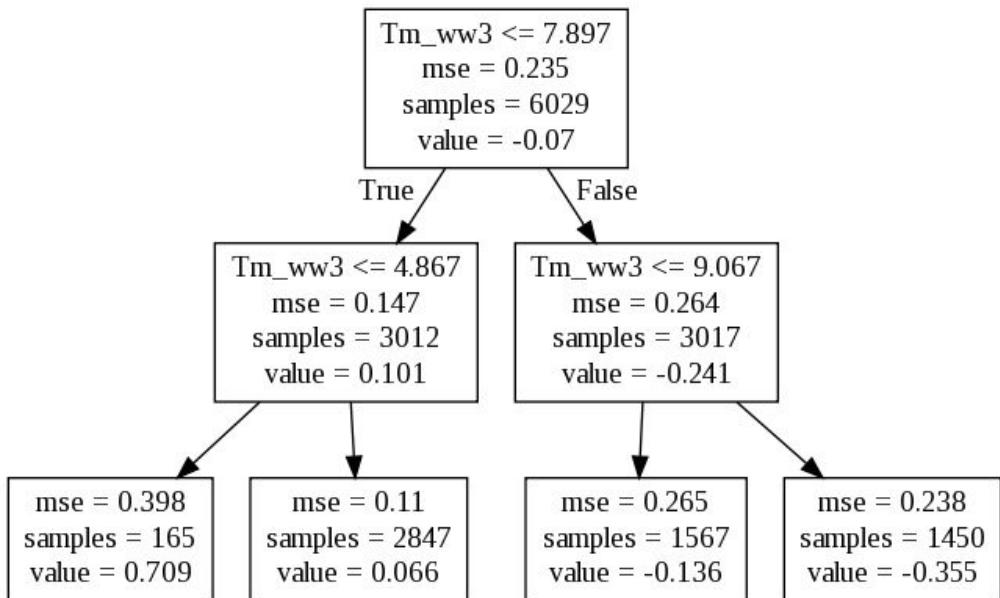
$$y(t) = H_{obs}(t) - H_{ww3}(t)$$



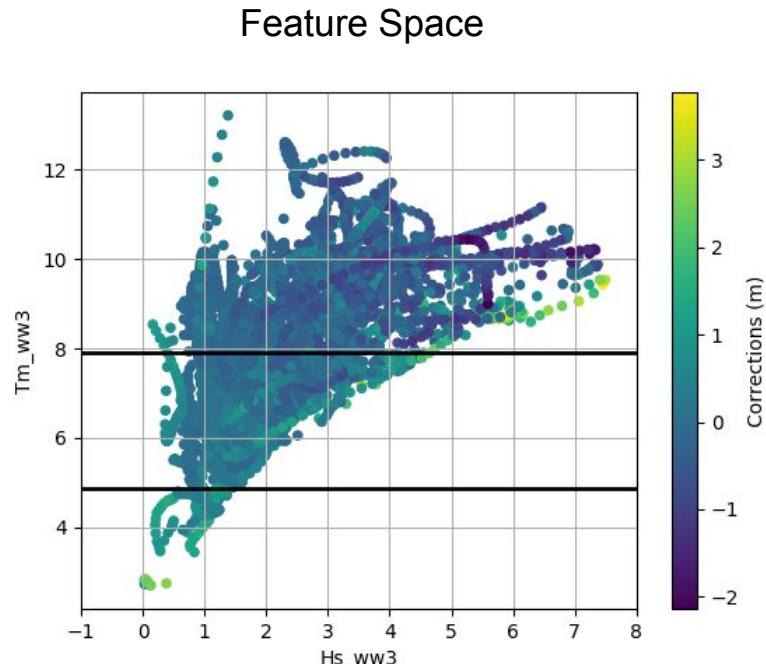
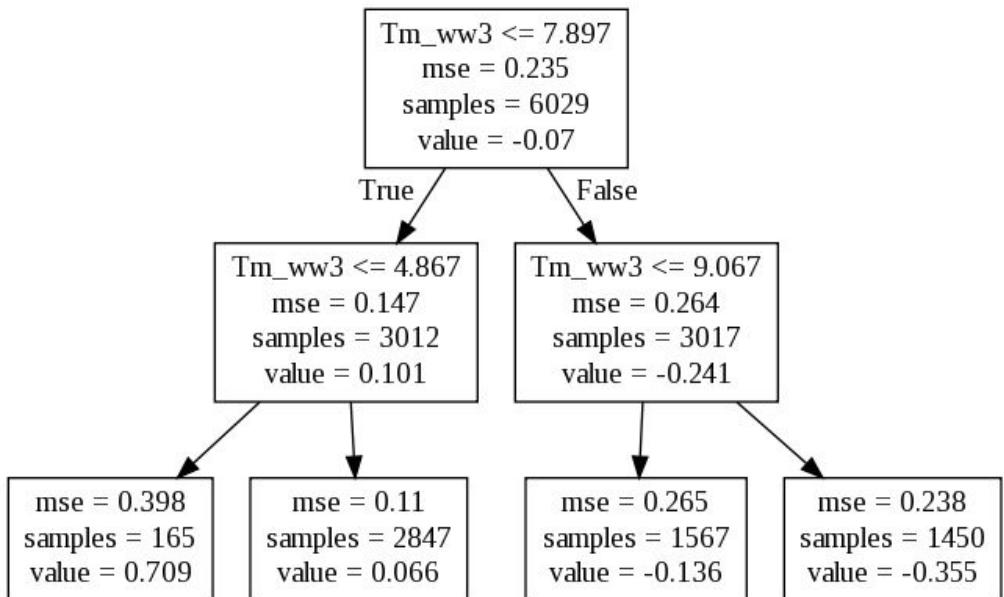
Decision Tree Logic



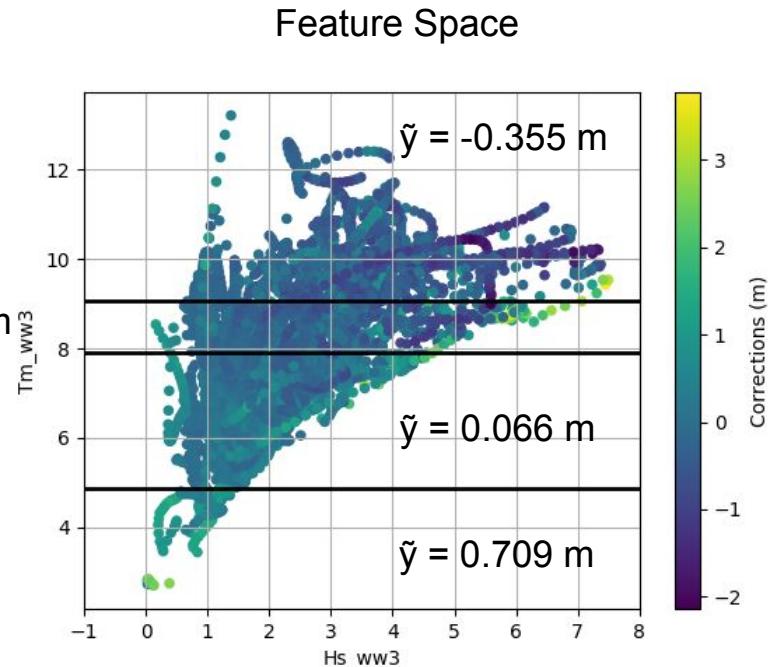
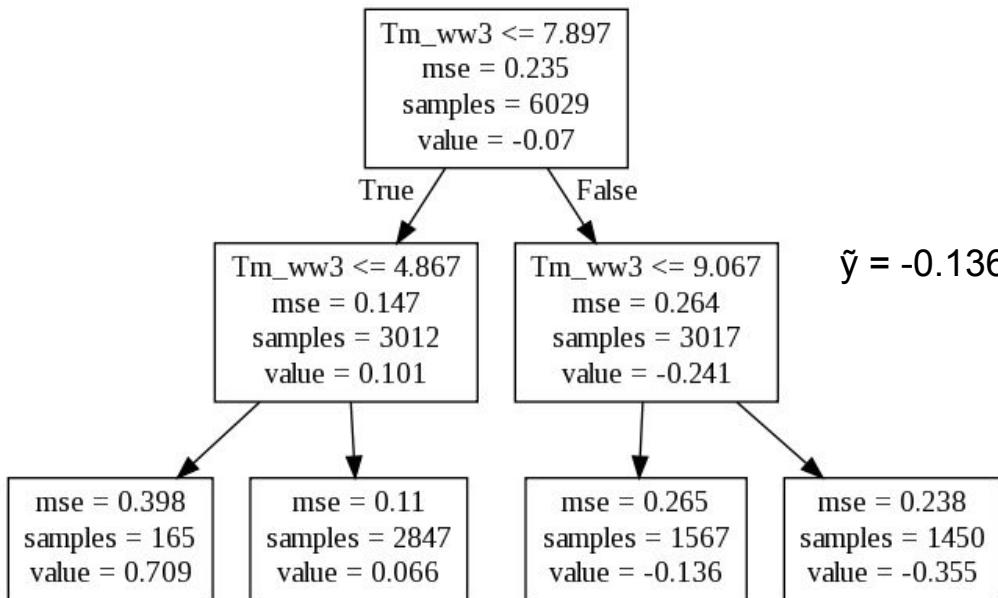
Decision Tree Logic



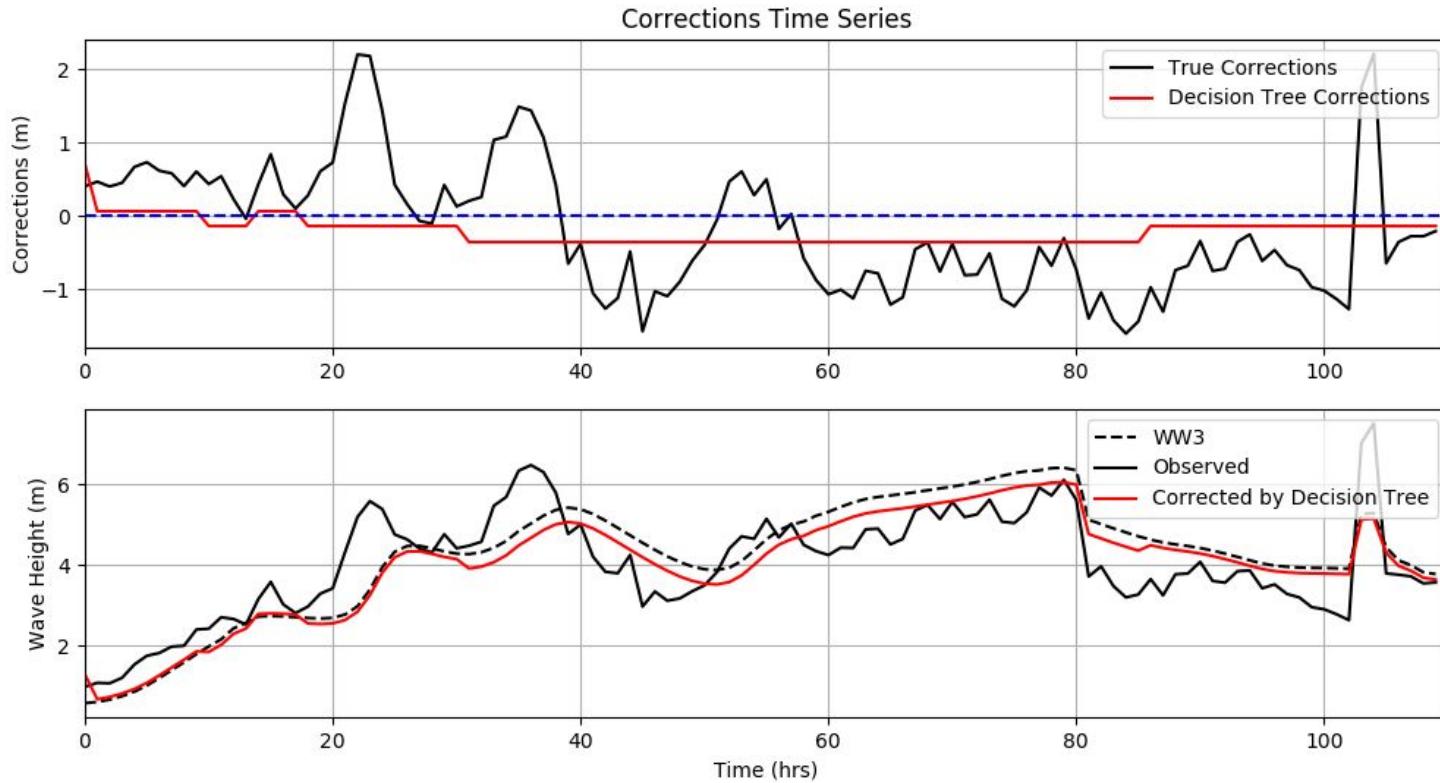
Decision Tree Logic



Decision Tree Logic



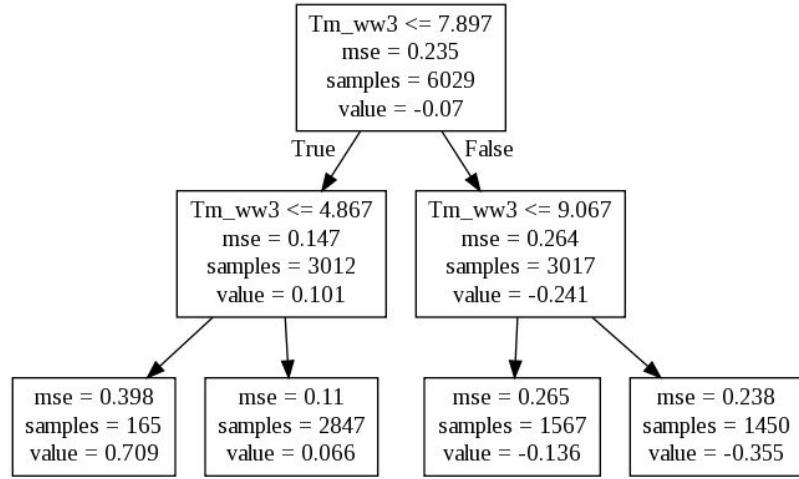
Decision Tree Output



Decision Tree On Wave Forecasts: A Complicated Case



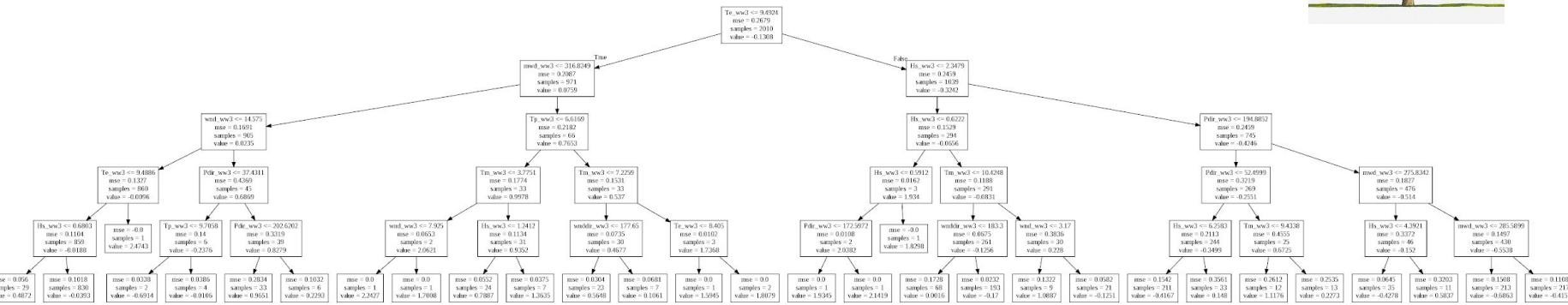
Random Forest Parameters: Tree Depth



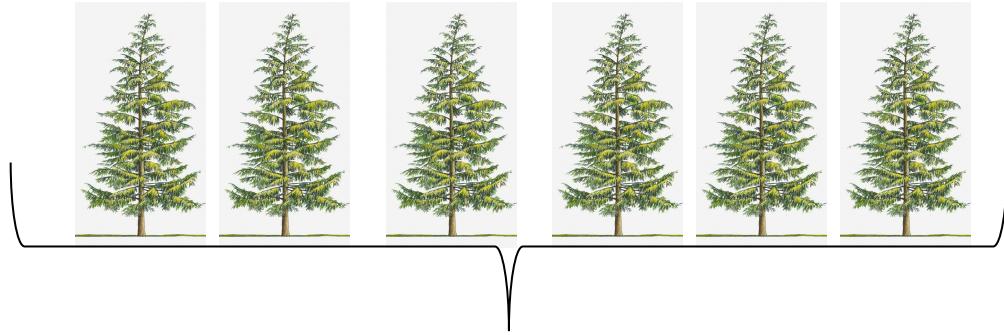
Tree depth: 2



Tree depth: 5



Bagged Decision Tree Parameters: Number of Trees



Final Prediction is average
of each tree's prediction

Bagging Method

Each tree is trained on a subset of the data.

Data points are picked uniformly and with replacement.

Determining Forest Parameters



VS



Decision Tree Variations for Different Problems

Ensemble Methods

- Bagging
- Boosting
- Random Forest

Prediction Methods

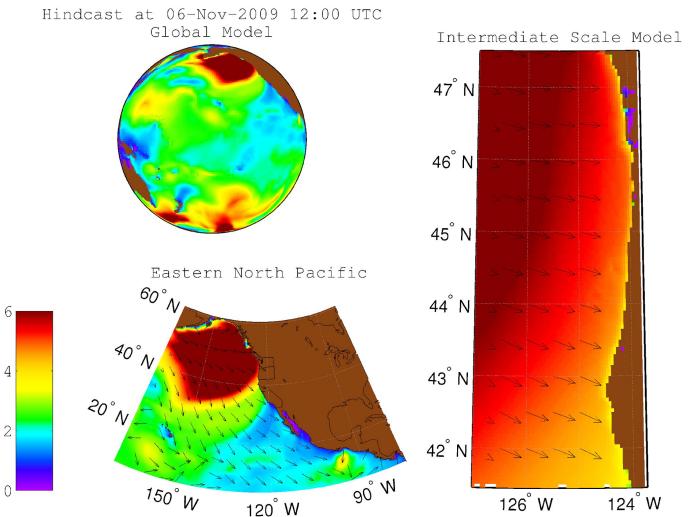
- Classification
- Regression

Objective Function (Regression)

- Mean Absolute Error
- Least Squares Error

Data

- Observations from NDBC Buoy 46050
- WW3 ST2 24-hour forecast horizon, years 2012-2015
- Forced by Global Forecast System (GFS) wind



Garcia-Medina et al, 2014

Experimental Set Up

Train

Summer (April 1-September 31)

- 10,158 points
- Years 2012-2014

Winter (Jan 1- March 31, October 1 - Dec. 31)

- N = 10, 092 points
- Years 2012-2014

Test

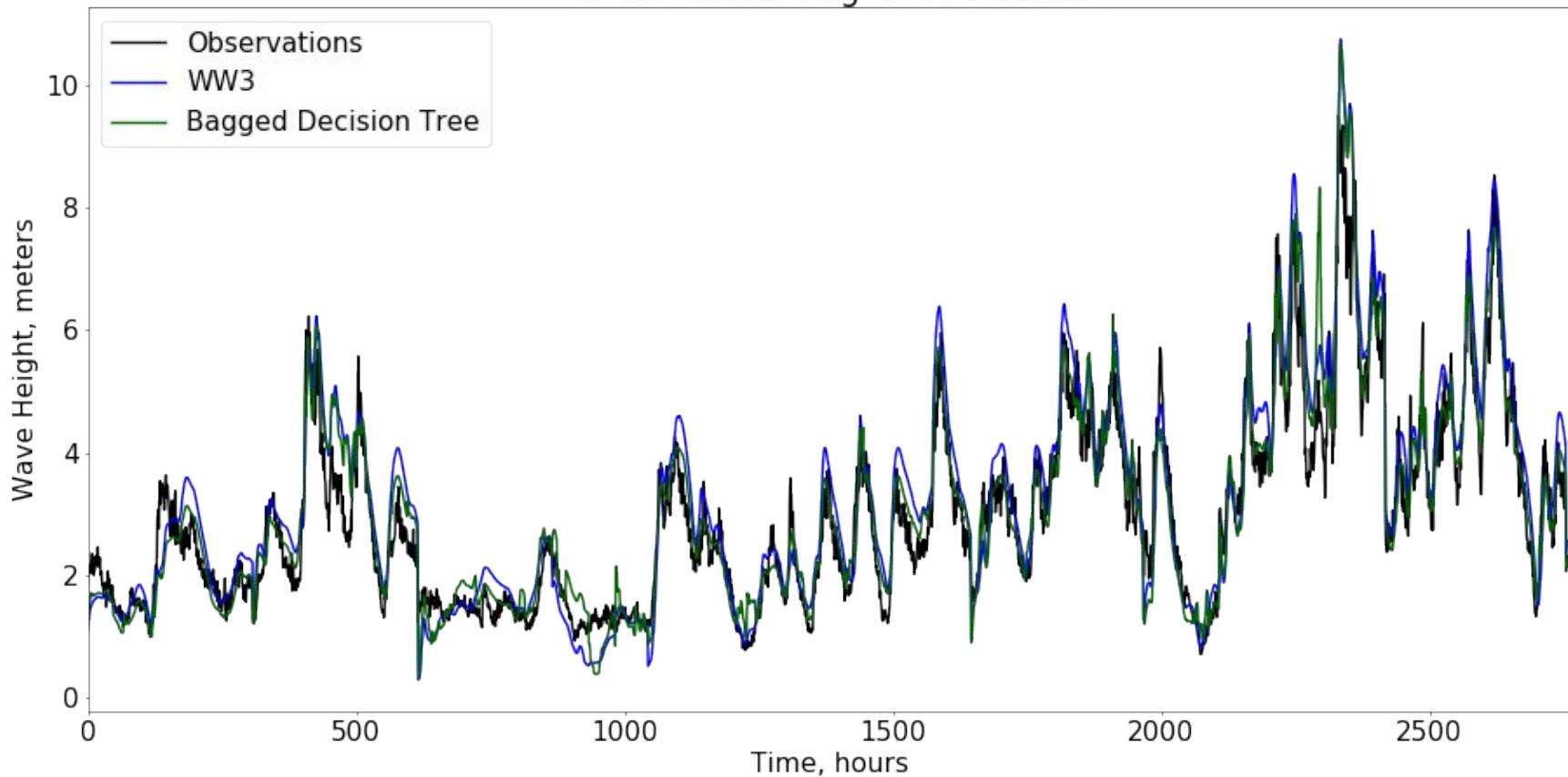
Summer 2015

- 2015
- 3,531 points

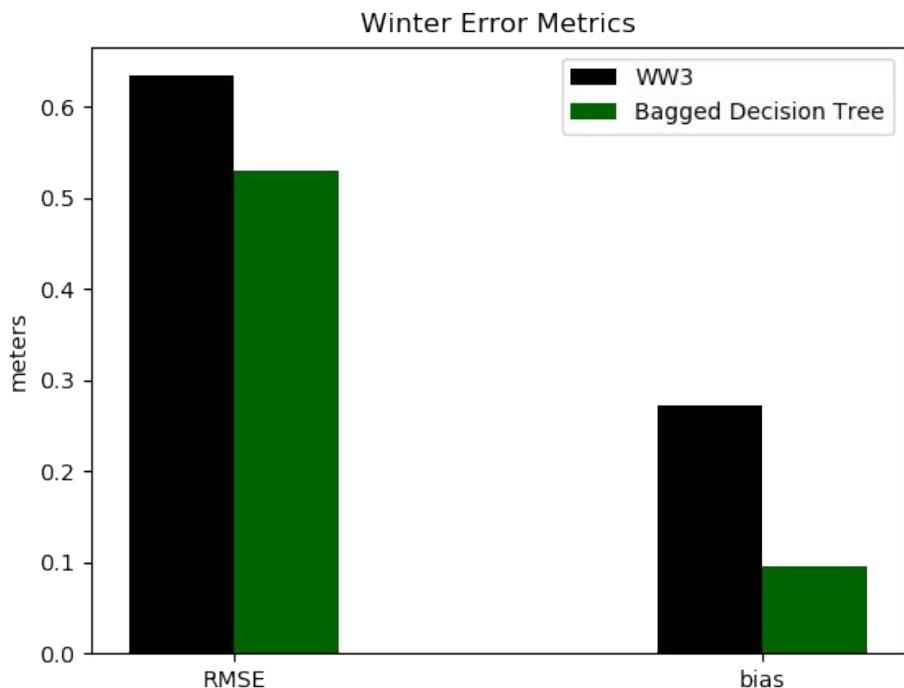
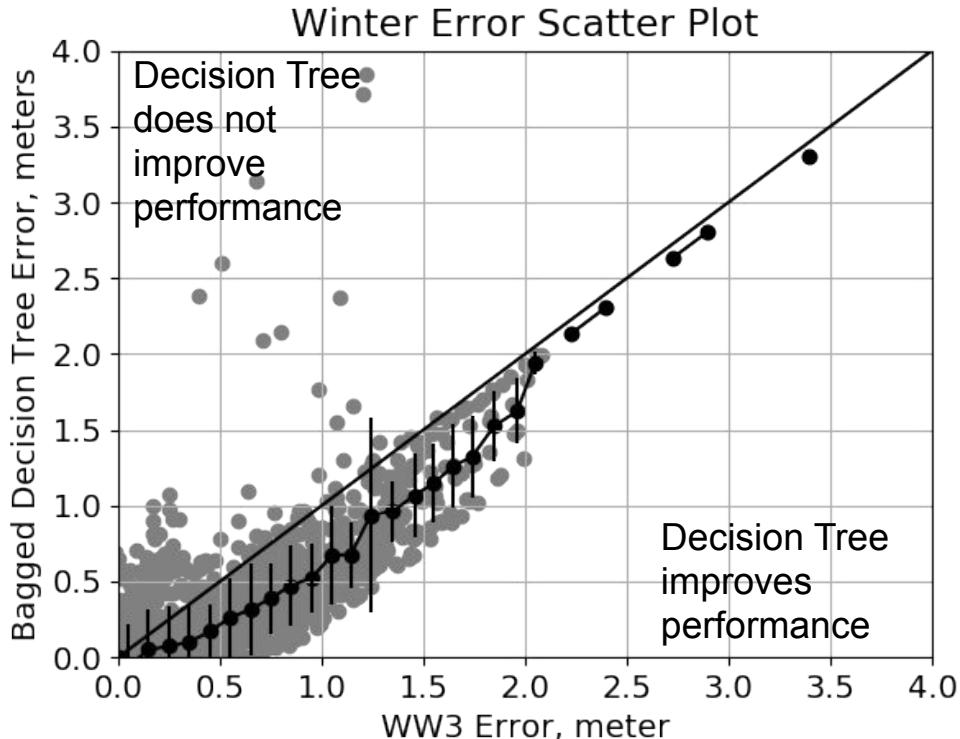
Winter 2015

- 2015
- 2,767 points

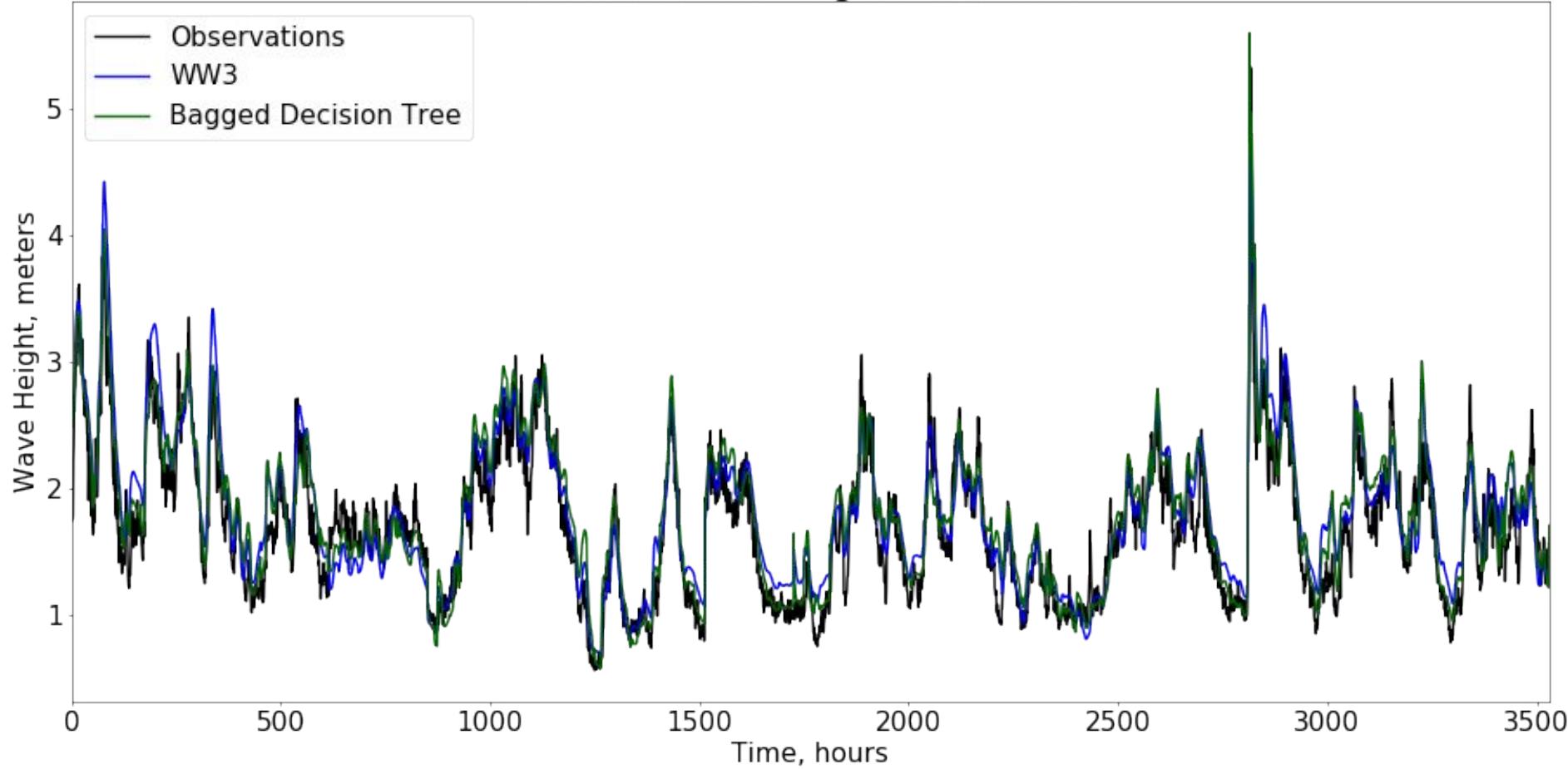
Winter Wave Height Time Series



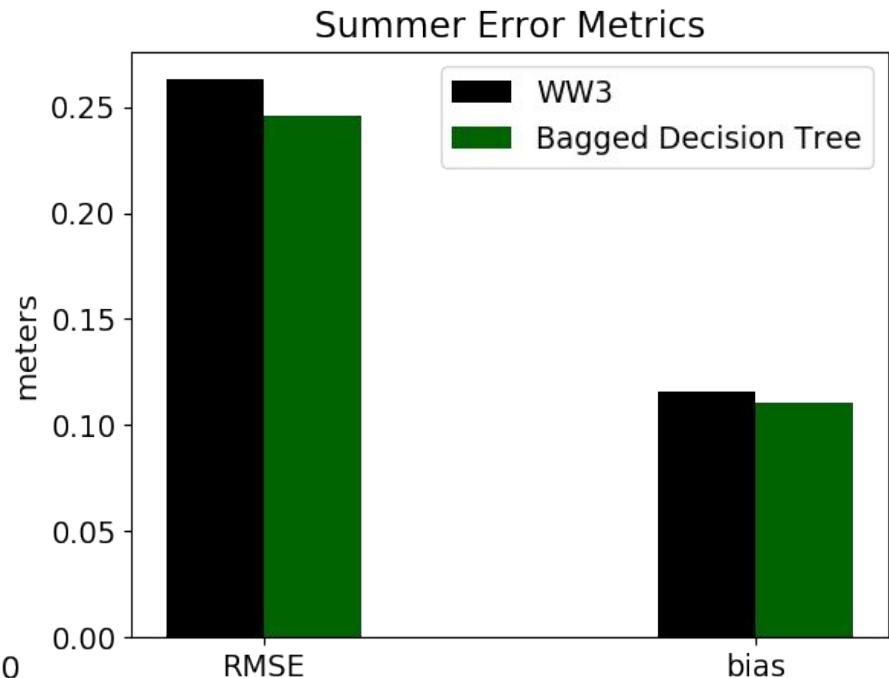
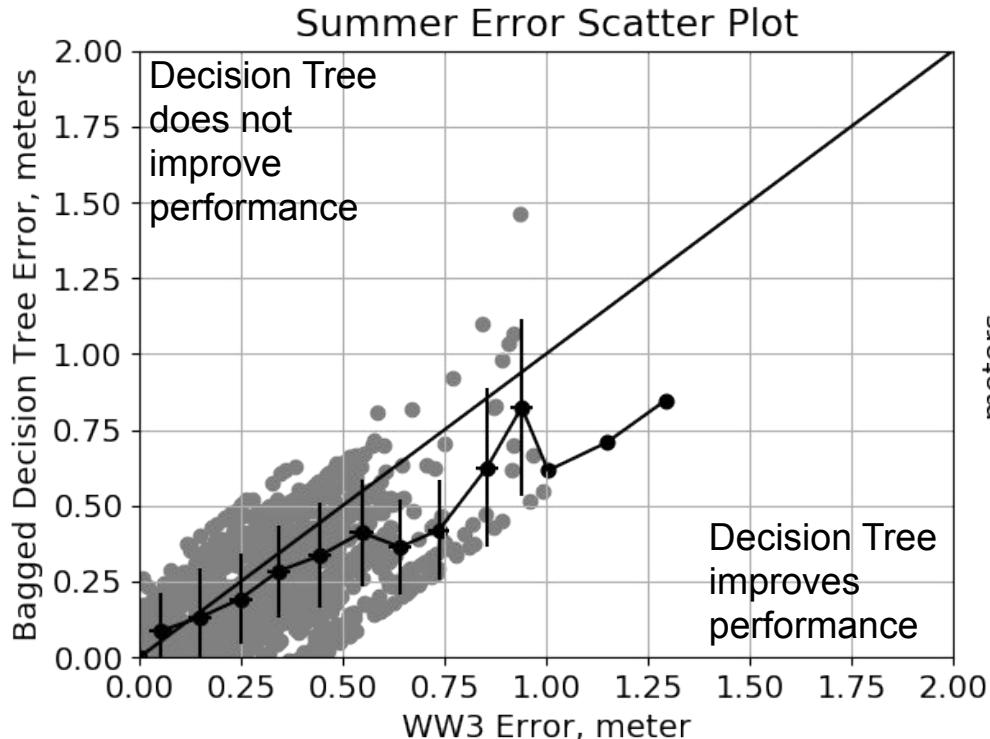
Results: Winter



Summer Wave Height Time Series

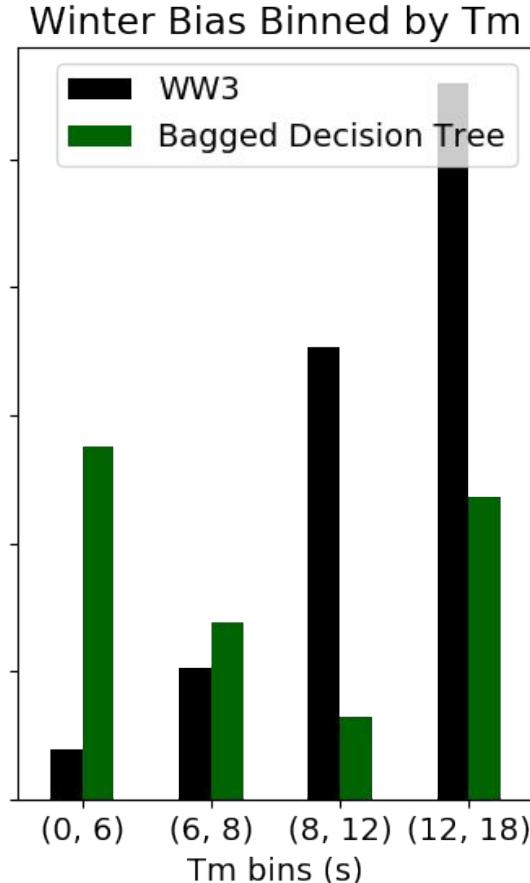
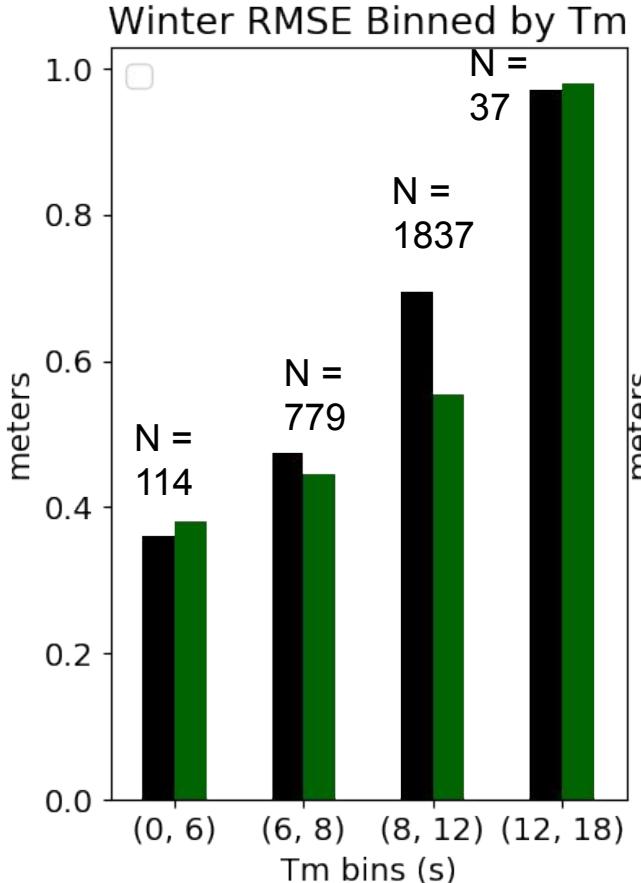


Results: Summer



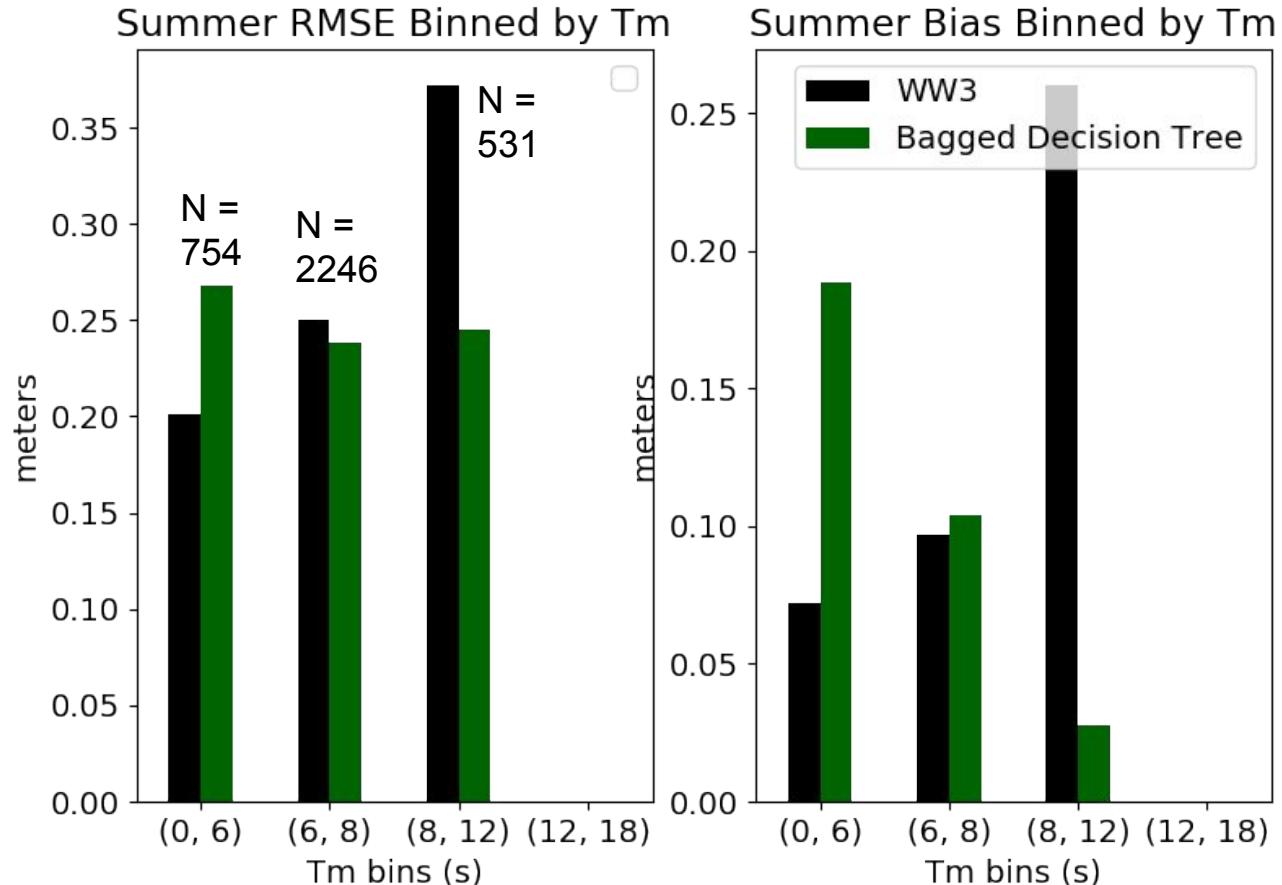
Results: Winter

- Decision Tree improves performance in mid-range periods (6-12s)
- Improves bias more than RMSE



Results: Winter

- Decision Tree improves performance in mid-range periods (6-12s)
- Improves bias more than RMSE



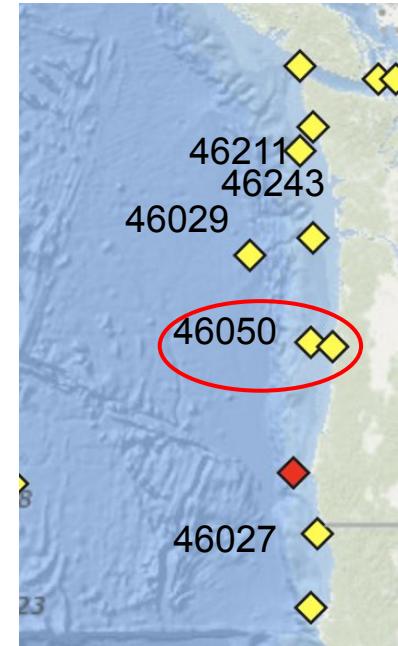
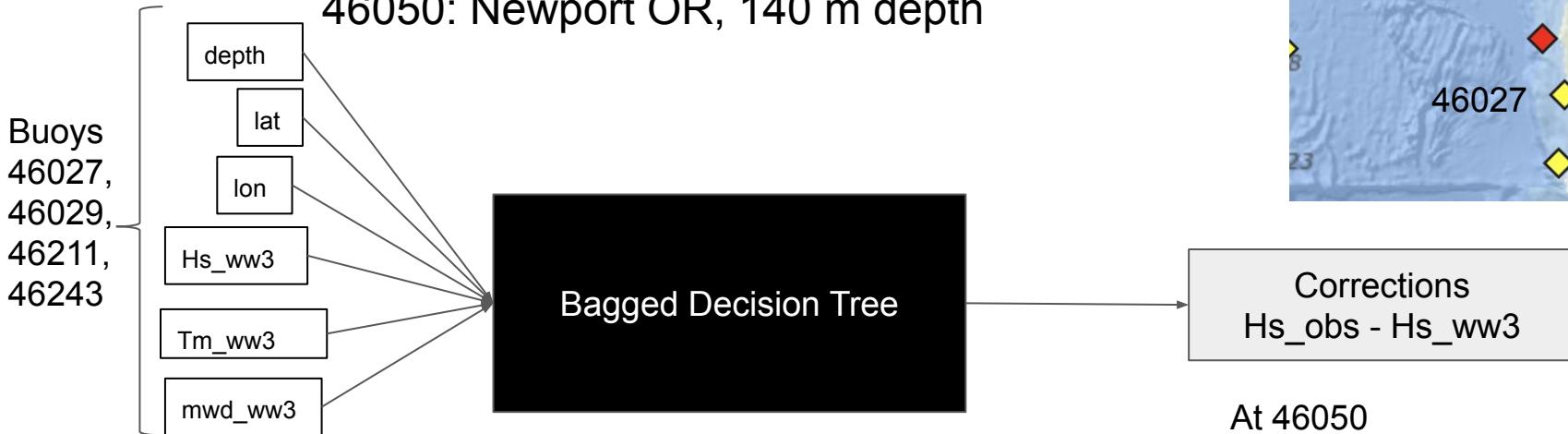
Can this be applied at different locations?

Train with data from other buoys:

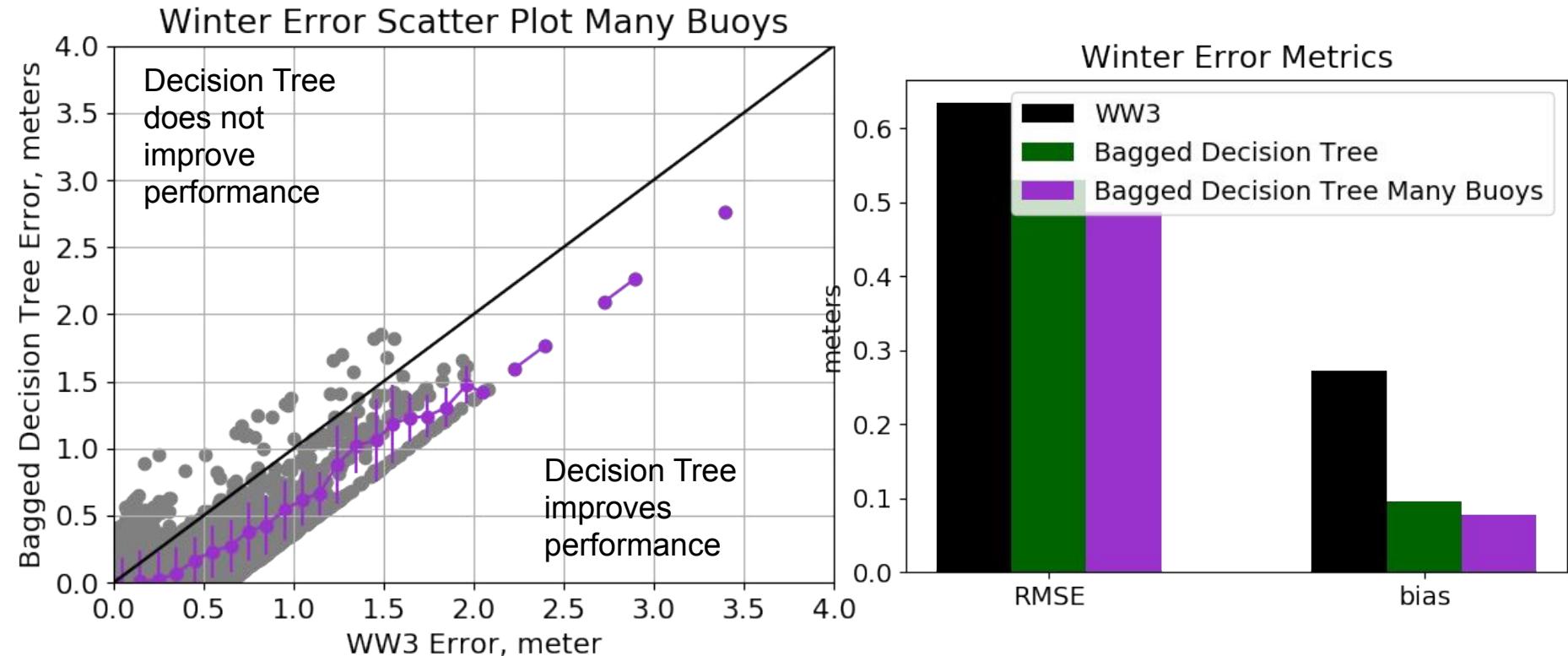
- 46027: Crescent City, CA, 46 m depth
- 46029: Columbia River, OR, 134 m depth
- 46211: Greys Harbor, WA, 40 m depth
- 46243: Clatspot Spit, OR, 24.4 m depth

Test on

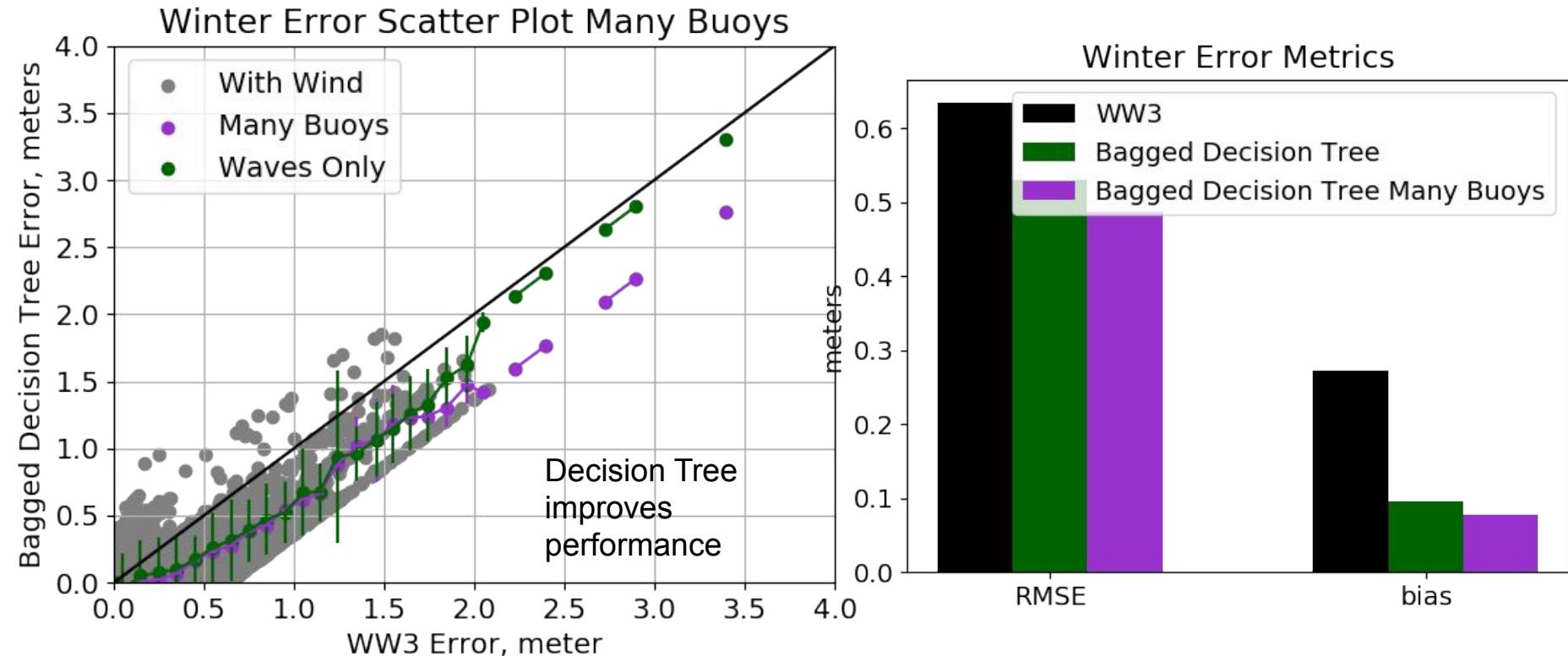
46050: Newport OR, 140 m depth



Results: Winter Many Buoys



Results: Winter Many Buoys

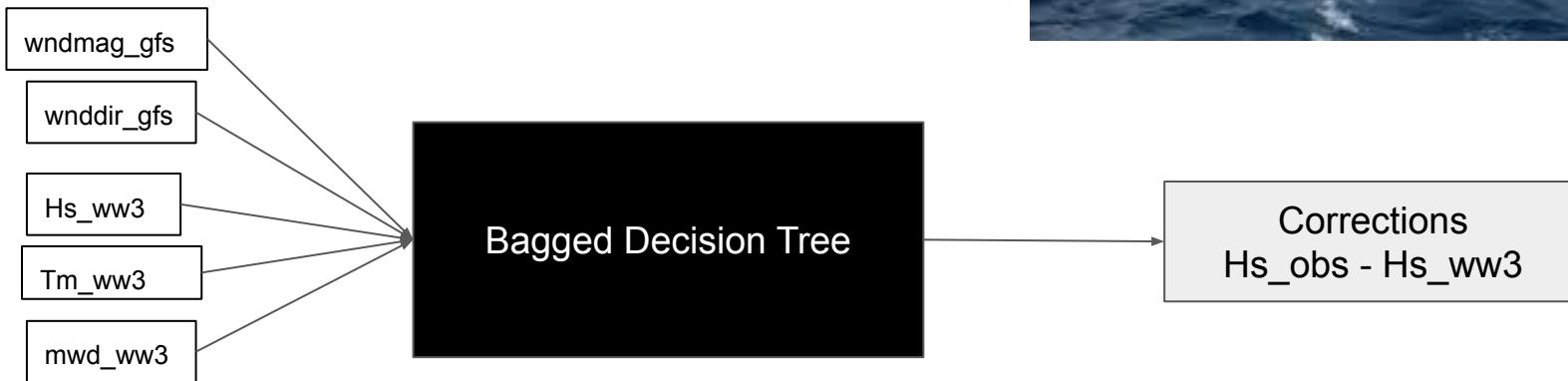


Can performance be further improved when more environmental information is included?

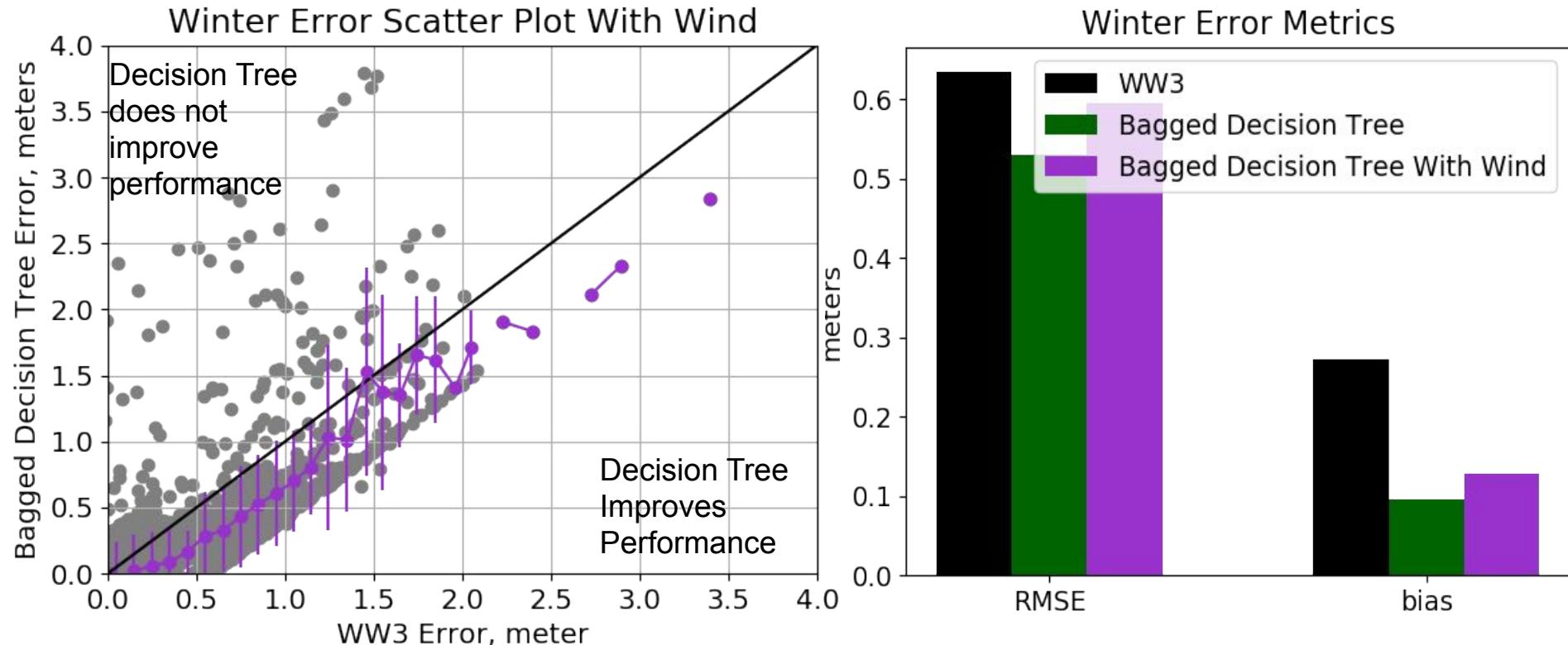
Information about wind forcing

- Wind direction
- Wind magnitude

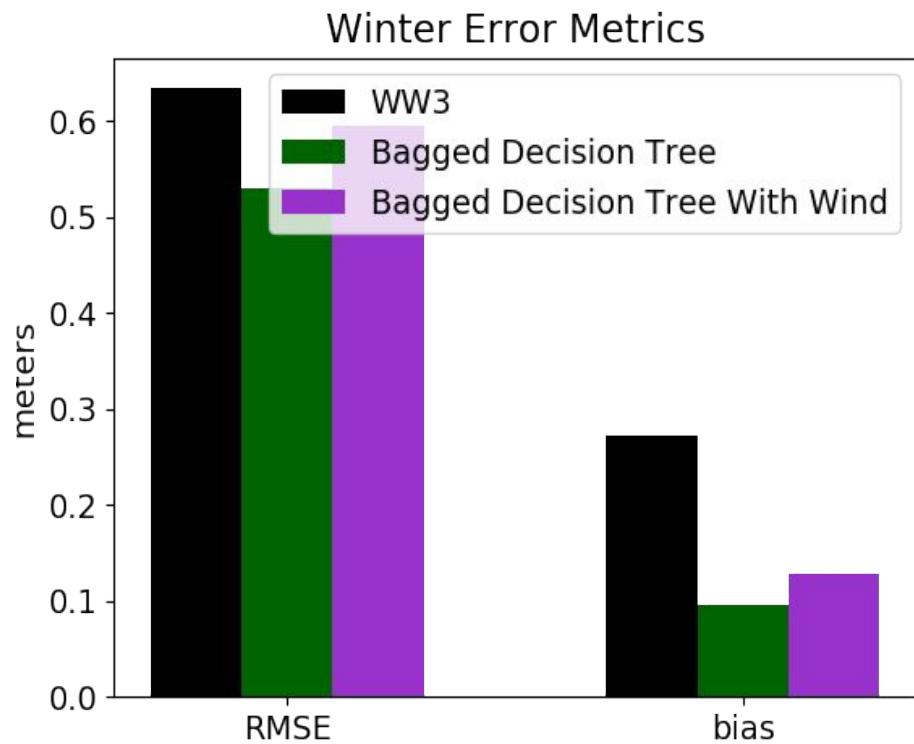
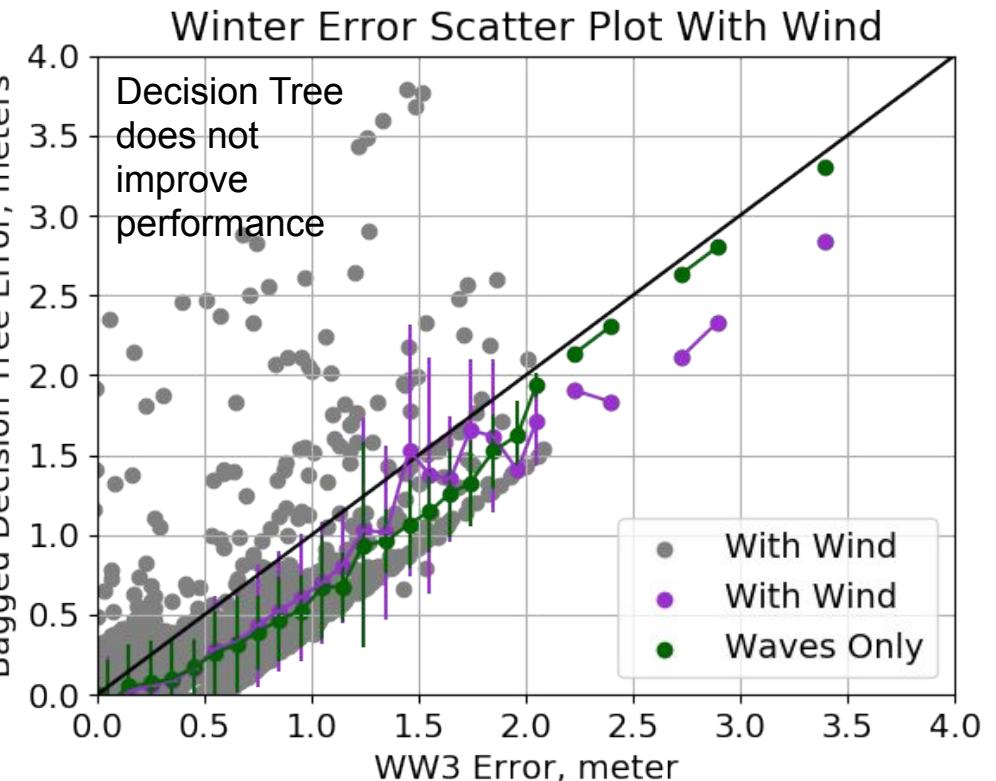
Data Source: Global Forecast System
(GFS) 24 hr forecast



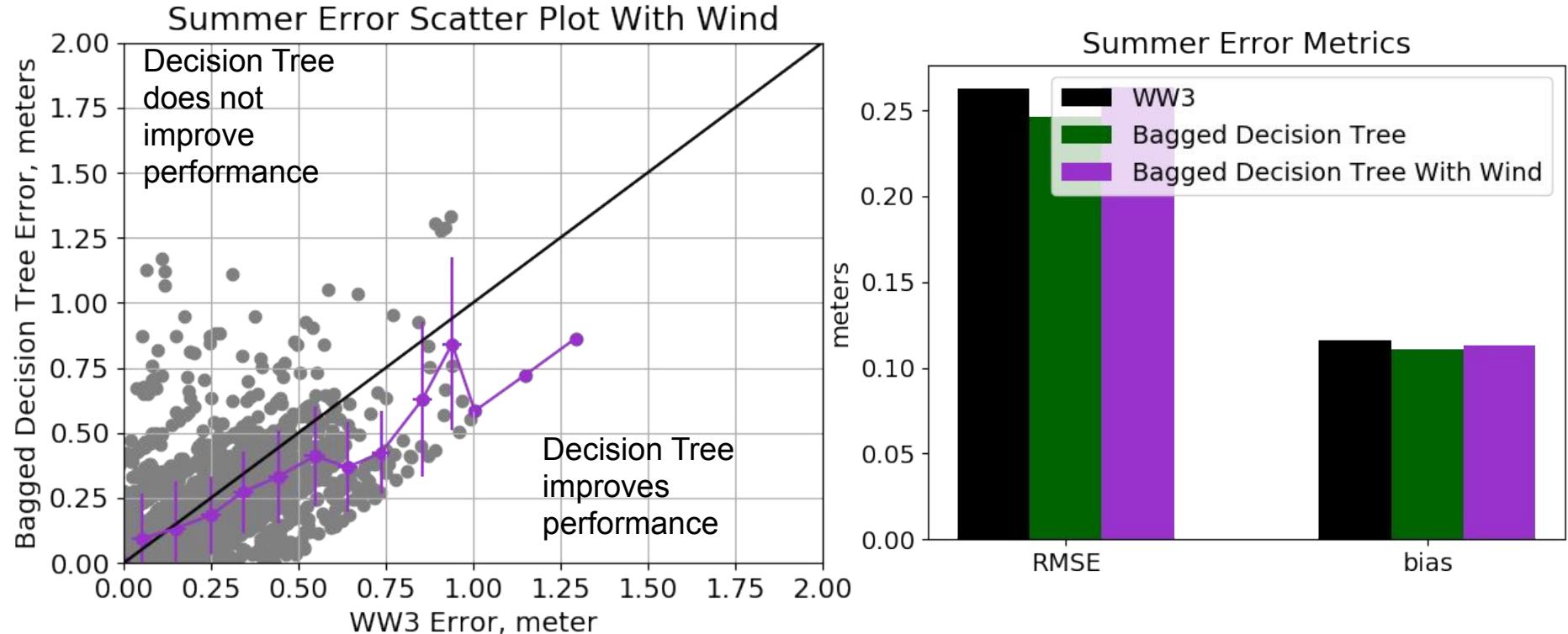
Results: Winter With Wind



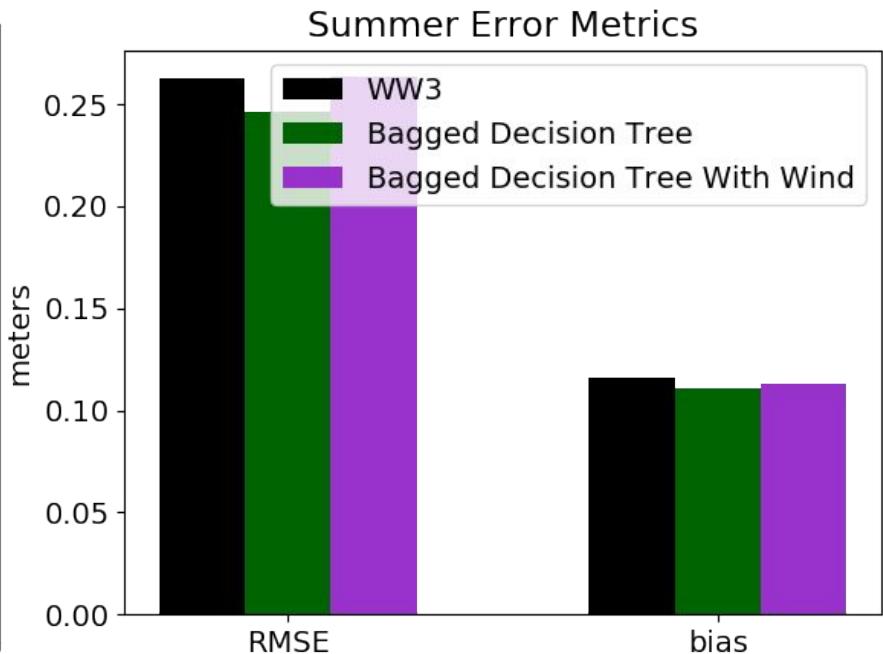
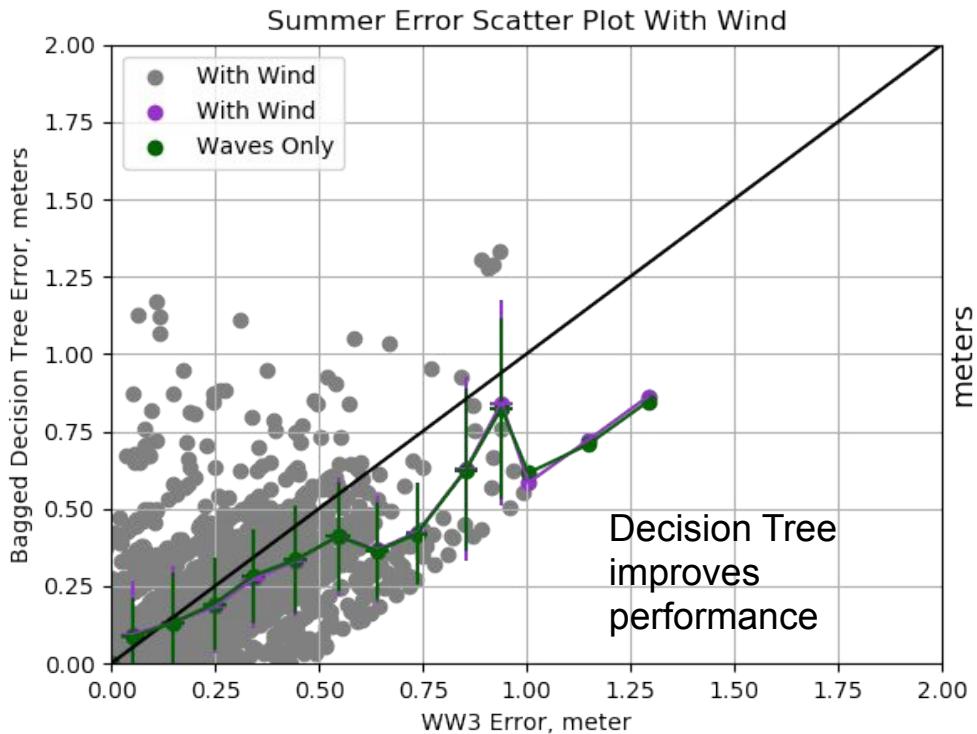
Results: Winter With Wind



Results: Summer With Wind

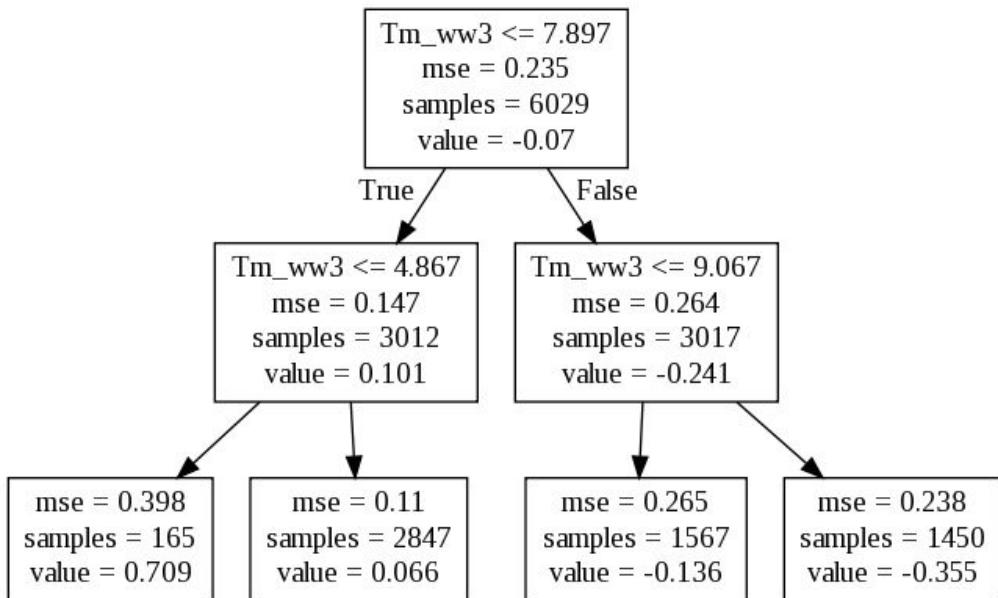
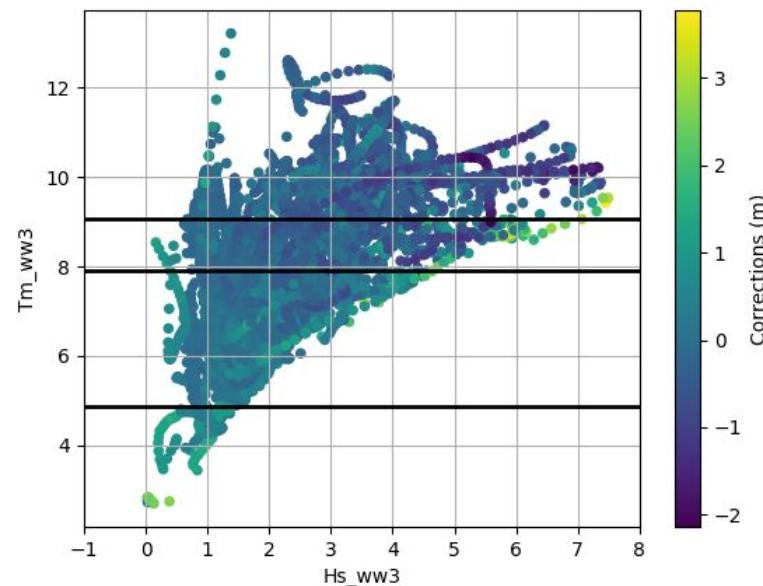


Results: Summer With Wind



Feature Importances

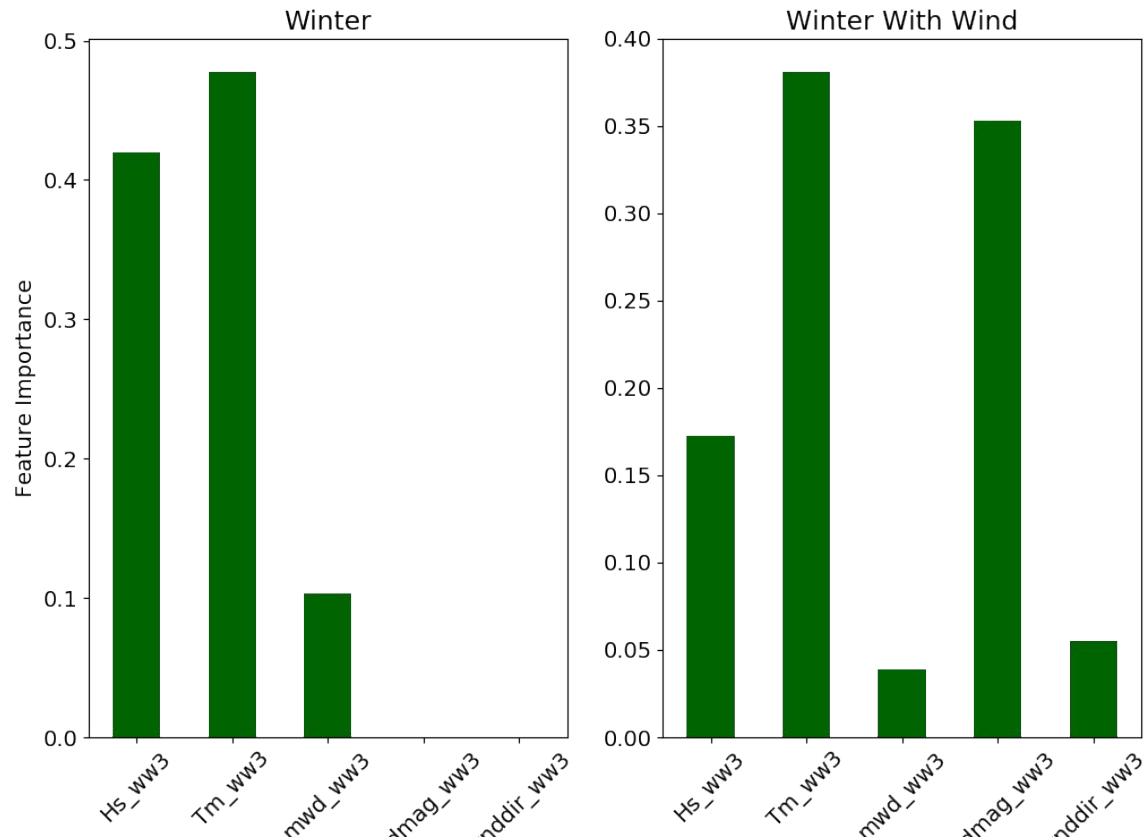
Feature Space



- In this case, DT only splits on Tm, not Hs
- Tm more important feature (feature importance of 1)

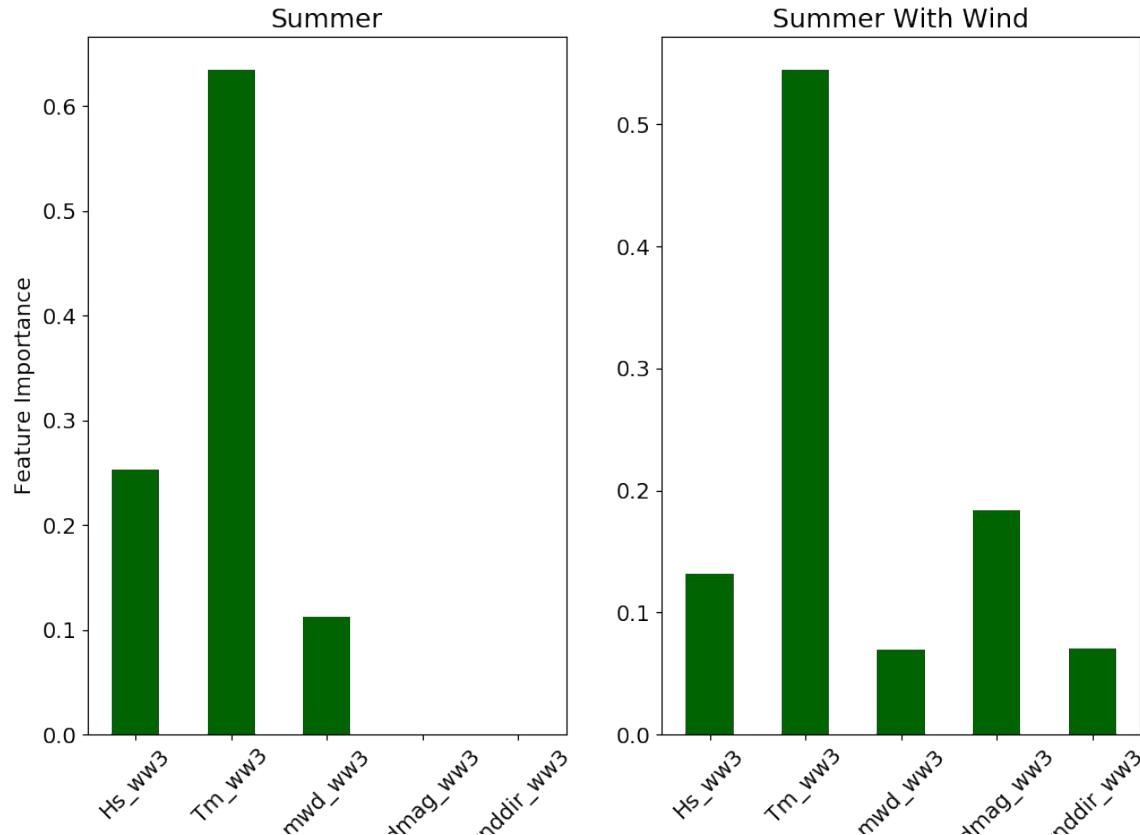
Feature Importances: Winter

- Mean period most important feature, Hs second
- When wind is added, wind magnitude becomes second most important feature

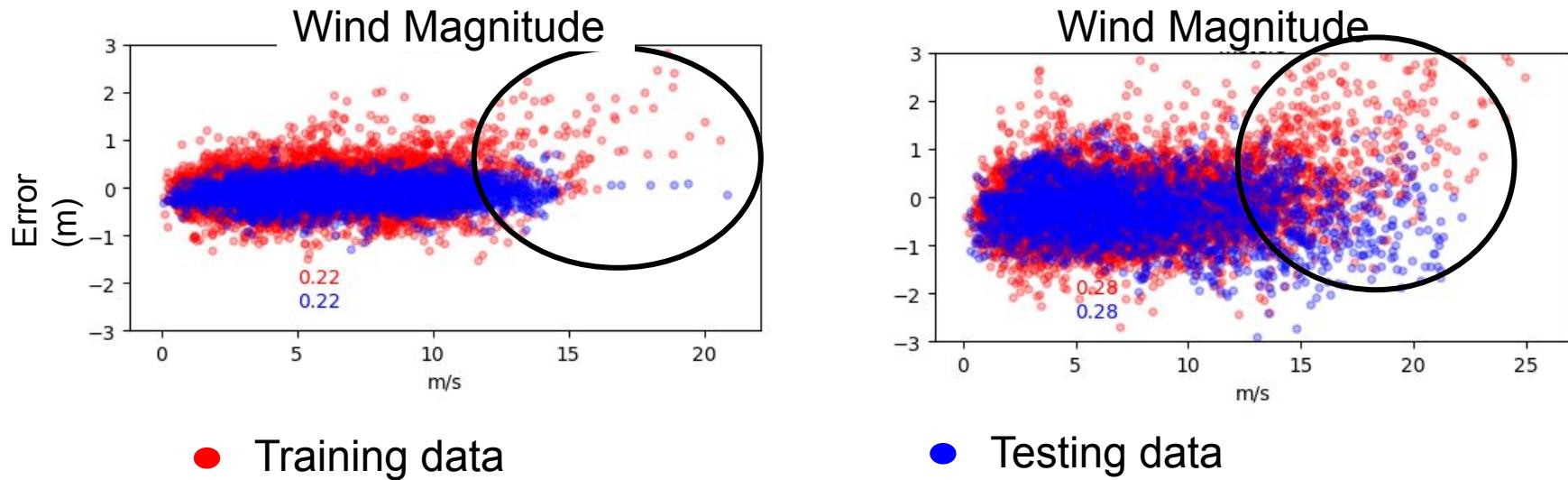


Feature Importances: Summer

- Mean period most important feature, Hs second
- When wind is added, wind magnitude becomes second most important feature



Training Data (2012-14) vs Testing Data (2015)



Hypothesis: The training wind magnitude values result in DT trained on environmental conditions that do not arise in test conditions

For best performance, provide a training set that is representative of the testing conditions - best to give it median conditions

Conclusions

General:

- Random forest able to find systematic error within numerical model output
 - Evidenced by decreased RMSE and bias of post processed time series
- Performs better for winter data than for summer
 - Hypothesis that greater variability in the data provides more opportunity for the decision tree to partition the data space

WW3 Systematic Errors:

- Mean period is the most important parameter for Random Forest to partition the data space
 - error correlation higher for longer period than shorter period

Performance:

- Train it on data which is representative of the testing data
- Don't extrapolate method out of training conditions