



O1 Context & Data Preparation

02 Multiple Regression

03 Classification Models

04 DataRobot & Conclusion



#### **Context**



- Trigger health issues
- Business problems



- Study the relationship
- The outcome optimized models help oversee and cut down

# **2** Data Wrangling

### Credible Data Source

- <u>https://www.kaggle.com/uciml/pm25-data-for-five-chinese-cities</u>
- "Assessing Beijing's PM2.5 pollution: severity, weather impact, APEC and winter heating", Liang, X., Zou, T., Guo, B., Li, S., Zhang, H., Zhang, S., Huang, H. and Chen, S. X. (2015)

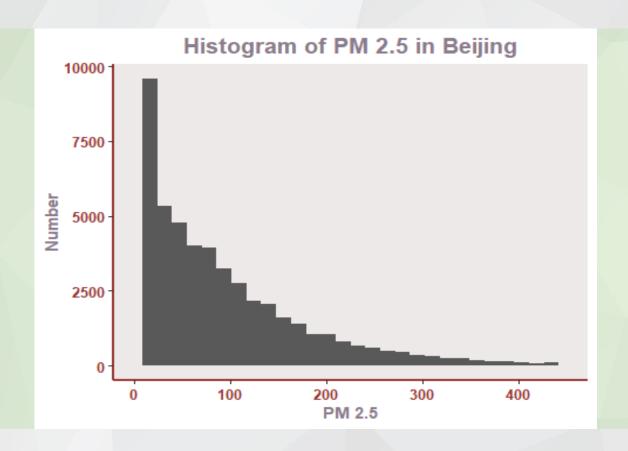
## Trifacta Data Wrangling

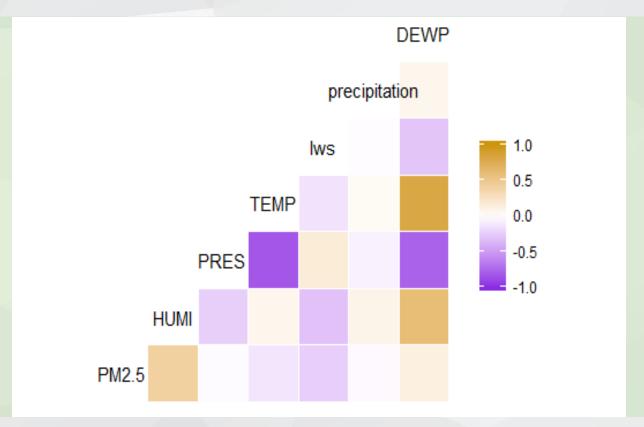
49,579 valid observations and 13 variables; Time period 2010-2015

- Drop missing values (5.7%)
- Create 4 dummy variables for wind direction (NE, NW, SE, SW)

Variables	Description
Year	
PM_US.Post	PM2.5 (ug/m^3)
HUMI	Humidity (%)
DEWP	Dew Point (Celsius Degree)
TEMP	Temperature (Celsius Degree)
PRES	Pressure (Pa)
seasonadj	Season
lws	Wind Speed (m/s)
precipitation	Precipitation (mm)
SW	Wind from southwest direction
SE	Wind from southeast direction
NW	Wind from northwest direction
NE	Wind from northeast direction

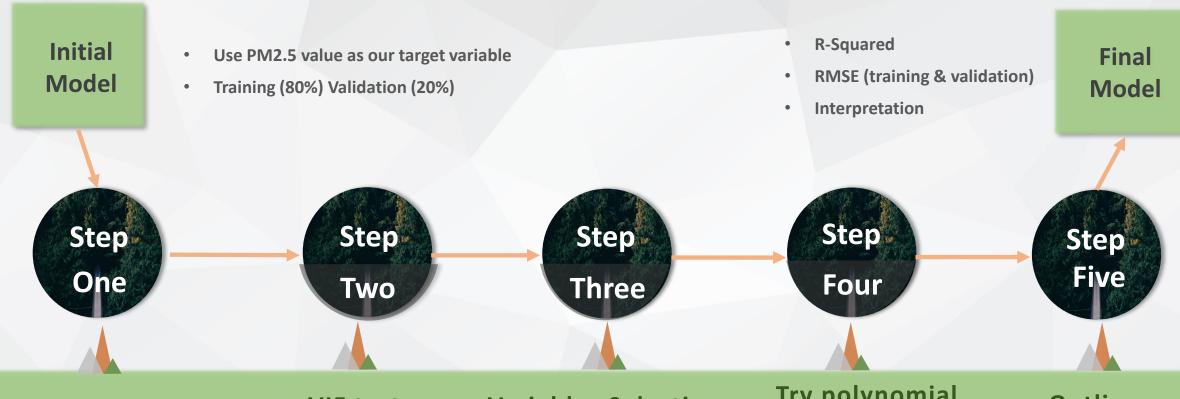
# Data Visualization







### **Multiple Regression-Modeling Process**



Add Logarithmic form

Log(PM 2.5)

#### VIF test

- Temperature > 5
- Dew Point > 5

#### **Variables Selection**

- forward selection. backward selection and exhaustive search
- 10 variables (delete temperature & dew point)

#### Try polynomial term in our model

- Air pressure^2
- Wind speed\*wind direction

#### **Outliers**

Delete 21 Obs (0.053%)



### Multiple Regression-Variable Selection

```
Coefficients:
                 Estimate Std. Error t value Pr(>|t|)
                                             <2e-16 ***
(Intercept) 1.984e+01 6.140e-01 32.311
HUMI
               2.022e-02 1.739e-04 116.279
                                             <2e-16 ***
seasonadjSummer -7.310e-01 1.285e-02 -56.869
                                             <2e-16 ***
seasonadjFall -2.820e-01 1.209e-02 -23.318
                                             <2e-16 ***
seasonadjWinter 3.547e-01 1.345e-02 26.375
                                             <2e-16 ***
               -2.964e-03 8.908e-05 -33.272
                                             <2e-16 ***
Iws
cbwdNW
               -1.262e-01 1.384e-02 -9.119
                                             <2e-16 ***
            6.476e-01 1.346e-02 48.127
                                             <2e-16 ***
cbwdSE
                                             <2e-16 ***
cbwdSW
         4.733e-01 1.428e-02 33.155
precipitation -5.723e-02 4.481e-03 -12.773
                                             <2e-16 ***
                                             <2e-16 ***
PRES
               -1.663e-02 6.038e-04 -27.536
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' '1
```

Residual standard error: 0.7797 on 39652 degrees of freedom Multiple R-squared: 0.4481, Adjusted R-squared: 0.448 F-statistic: 3219 on 10 and 39652 DF, p-value: < 2.2e-16

Forward selection, backward selection:

10 variables: humidity, wind speed, precipitation, air pressure, three season dummies, three wind direction dummies.

Exhaustive search:

set maximum variables number: 10, (still add 12 variables into the model.)
Delete temperature and dew point

## Multiple Regression-Polynomial Terms



Wind speed\* SW, Wind speed\* SE, Wind speed\* NW

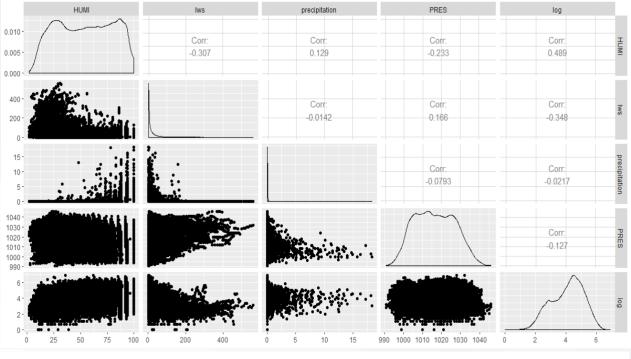
Hourly Wind Speed\*combined wind direction (categorical variable).

Wind from some specific directions may have more influence on PM2.5 in Beijing

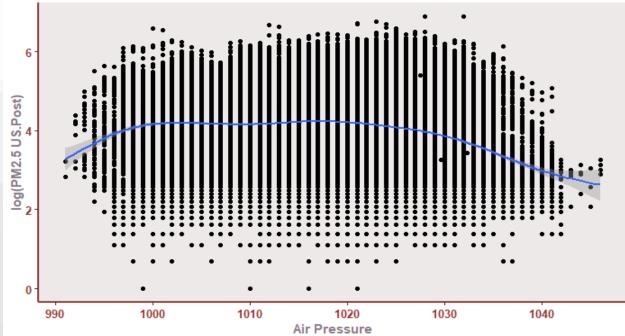


#### Pressure ^2

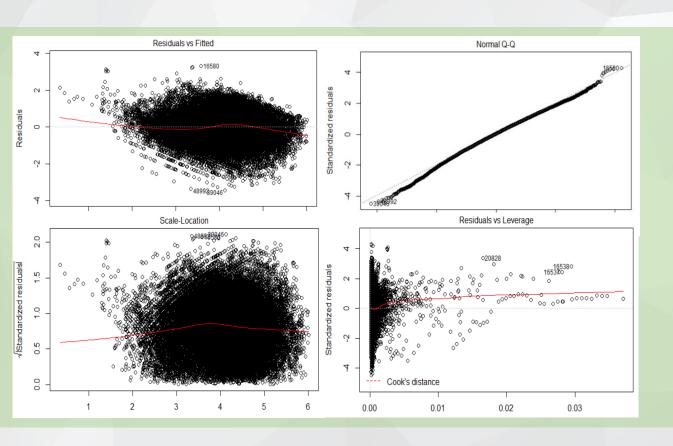
Plots between numeric variables and log(PM 2.5) log(PM 2.5) - Air Pressure - quadratic function log(PM 2.5) - Air Pressure^2 – linear relationship

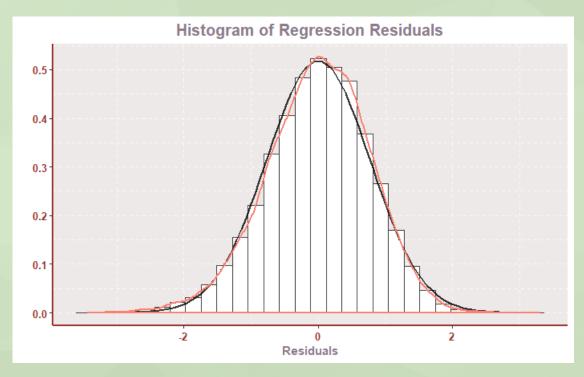


Beijing: The Relationship between Pressure and log(PM2.5\_US.Post)



## Multiple Regression-Outliers





# 5

Coefficients:

#### Multiple Regression-Interpretation

```
Estimate Std. Error t value Pr(>|t|)
(Intercept)
                           1.744e-04 113.439
HUMI
seasonadjSummer -6.368e-01
                           1.341e-02 -47.481
seasonadiFall
                           1.212e-02 -24.529
seasonadjWinter 3.930e-01
                           1.345e-02
                           3.640e-03 -16.662 < 2e-16
               -6.065e-02
Iws
cbwdNW
               -2.532e-01
                           1.519e-02 -16.667 < 2e-16
cbwdSE
                           1.555e-02
                                      31.901 < 2e-16
                                      22,101
cbwdSW
               -1.176e-01
                           7.111e-03 -16.535
precipitation
I(PRES^2)
               -7.498e-04
                           3.838e-05 -19.537
               1.510e+00
                           7.812e-02
                                      19.334
PRES
I(Iws * SW) 1.313e-02
                           1.944e-03
                                       6.751 1.49e-11
I(Iws * SE)
              1.581e-02
                           8.529e-04
                                      18.536
I(Iws * NW)
                1.445e-02
                          7.929e-04
                                     18.221
               0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Signif. codes:
Residual standard error: 0.7697 on 39627 degrees of freedom
                               Adjusted R-squared: 0.4614
Multiple R-squared: 0.4616,
F-statistic: 2427 on 14 and 39627 DF, p-value: < 2.2e-16
                   ME
                           RMSE
Test set -6.565243e-17 0.7695204 0.6098349
```

Test set -0.0006729632 0.7770129 0.616476 -Inf

Adjusted R^2 :0.4614 RMSE(Training): 0.7695 RMSE(Validation): 0.7770

#### **Humidity:**

coefficient = 0.01978 > 0.

When humidity increases by 1%, PM2.5 increases by 1.978%.

#### Wind direction:

(compared with wind from northeast)

#### PRES and PRES^2:

Coefficient(PRES) is positive, Coefficient(PRES^2) is negative. (under extreme situations, PM2.5 tends to have a negative relationship with pressure.)



#### **Target Variable Justification**

#### **Target variable: Air Condition**

If PM2.5 is less than or equal to  $50 \mu g/m^3$ , we consider it to be a 'Good' air condition, which is only unhealthy for sensitive groups.

If PM2.5 is more than 50µg/m³, we consider it to be a 'Bad' air condition, which is unhealthy for everyone.

AQI Category	Index Values	Revised Breakpoints (μg/m³, 24-hour average)
Good	0 - 50	0.0 <b>- 12.0</b>
Moderate	51 - 100	12.1 – <b>35.4</b>
Unhealthy for Sensitive Groups	101 – 150	35.5 – <b>55.4</b>
Unhealthy	151 – 200	55.5 – 150.4
Very Unhealthy	201 – 300	150.5 – 250.4
Hazardous	301 – 400	250.5 - 350.4
Hazardous	401 – 500	350.5 – 500

Image source: <a href="https://medium.com/mongolian-data-stories/air-pollution-part-2-f9f4da33a1bd">https://medium.com/mongolian-data-stories/air-pollution-part-2-f9f4da33a1bd</a>

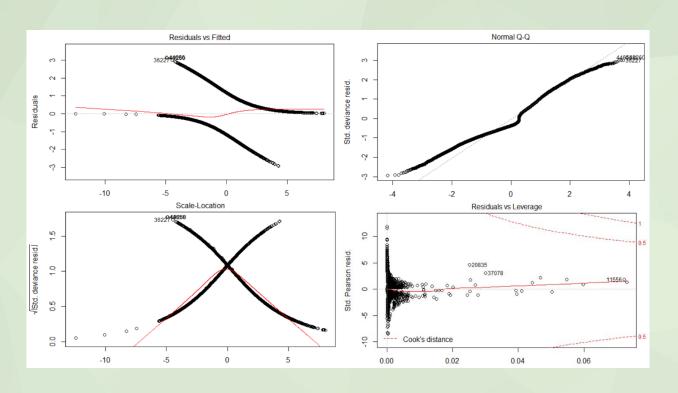


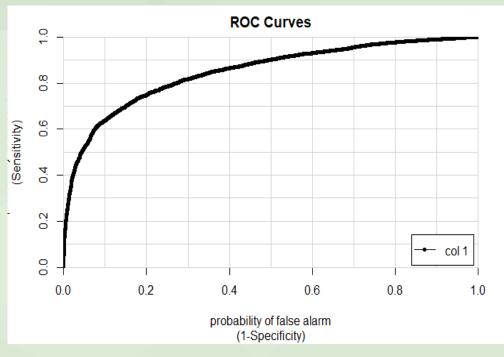
### **Logistic Regression-Modeling Process**





### **Logistic Regression-Final Model**





Test Set		Reference	
		Bad	Good
Prediction	Bad	8375	2235
	Good	737	3527

Accuracy:0.8002 (0.7969)Sensitivity:0.9191 (0.9156)F1:0.8493 (0.8450)AUC:0.8499



### **Logistic Regression-Coefficient Interpretation**

#### Odds Ratios of Logistic regression

Covariates	Logit coef.	Odds Ratio
Humidity	-0.050	0.951
Air pressure	-4.227	0.015
Air pressure^2	0.002	1.002
North west	0.280	1.323
South east	-1.830	0.160
South west	-1.162	0.313
Wind speed	0.010	1.010
Precipitation	0.953	2.594
Precipitation^2	-0.077	0.926
Summer	1.641	5.160
Fall	1.014	2.756
Winter	-0.775	0.461



**Positive Correlation:** Humidity, South wind, Winter.

**Negative Correlation:** North wind, Wind speed, Summer, Fall.

Air Pressure: Always negative.

**Precipitation:** First positive and then

becomes negative.

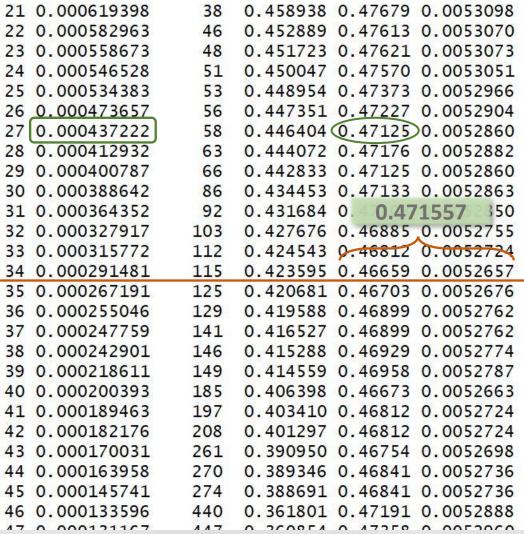
# 5

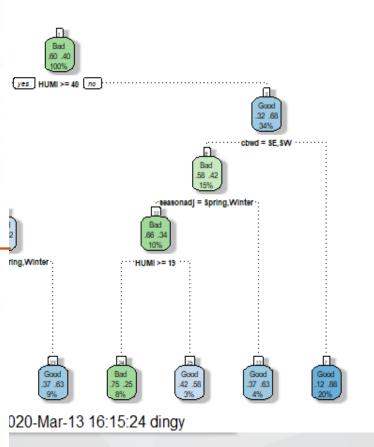
#### **Decision Tree-Pruning Process**

**Initial Model** 



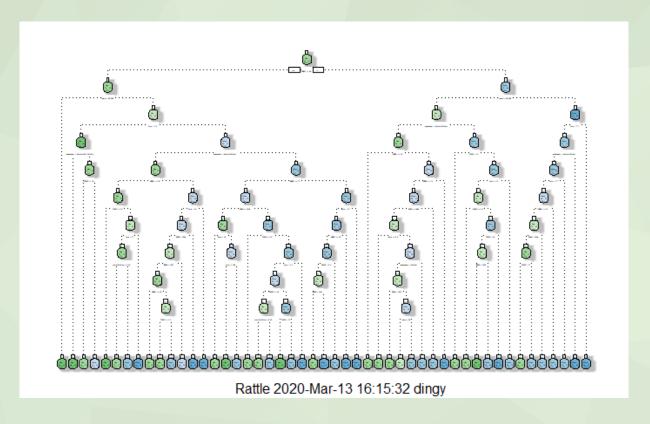
- Training (70%)Validation (30%)
- 'Bad' weather as positive class

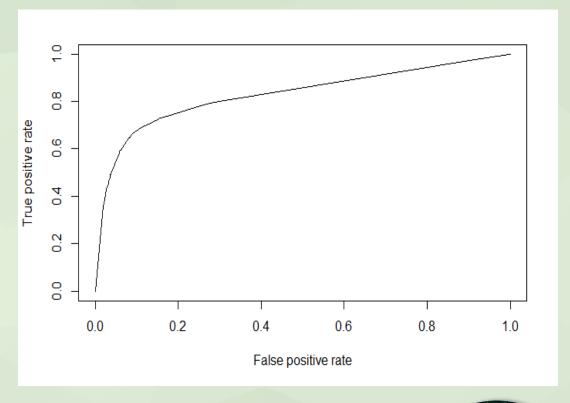






### **Decision Tree-Final Model**

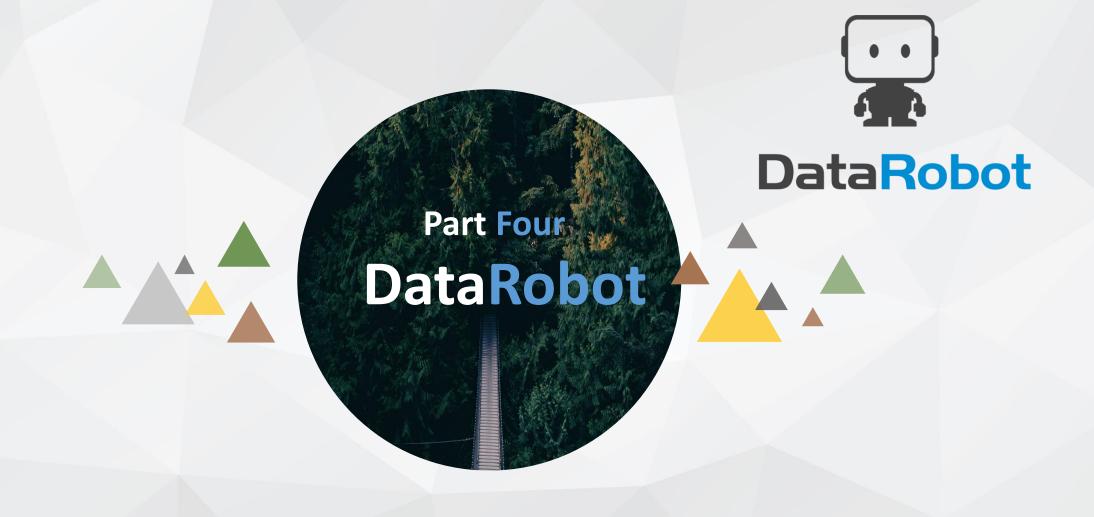




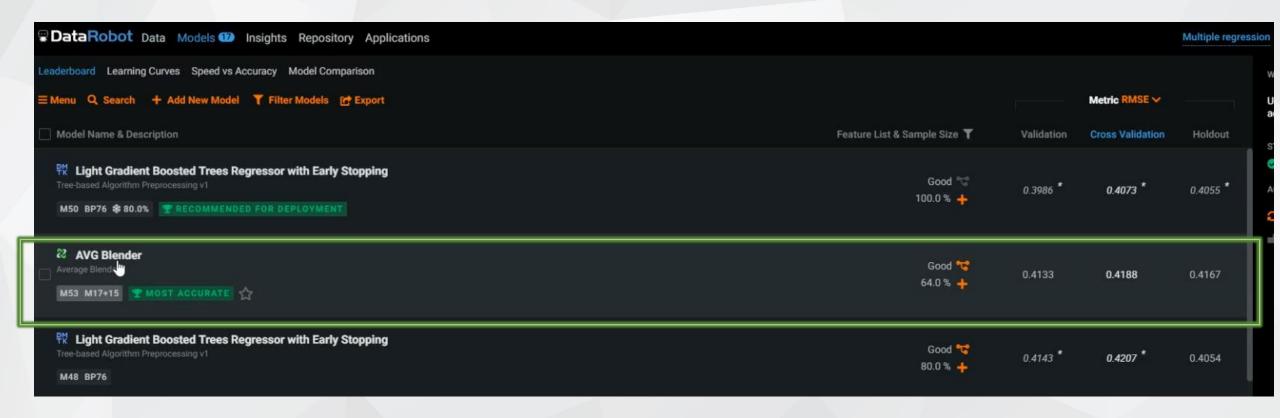
Test Set		Reference	
		Bad	Good
Prediction	Bad	8246	1890
	Good	866	3872

Accuracy: 0.8147 (0.8214)
Sensitivity: 0.9050 (0.9093)
F1: 0.8570 (0.8602)
AUC: 0.8290





### **DataRobot-Comparison of Regression Models**

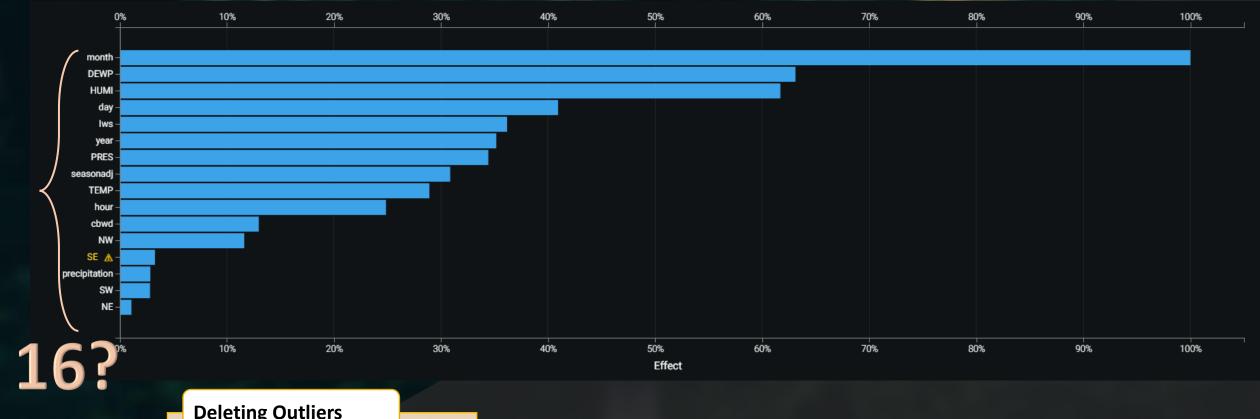


The AVG Blender has a RMSE of **0.4167**Our best handmake model has a RMSE of **0.78** 

Something Important!



### **DataRobot-Feature Importance**



**Deleting Outliers** 

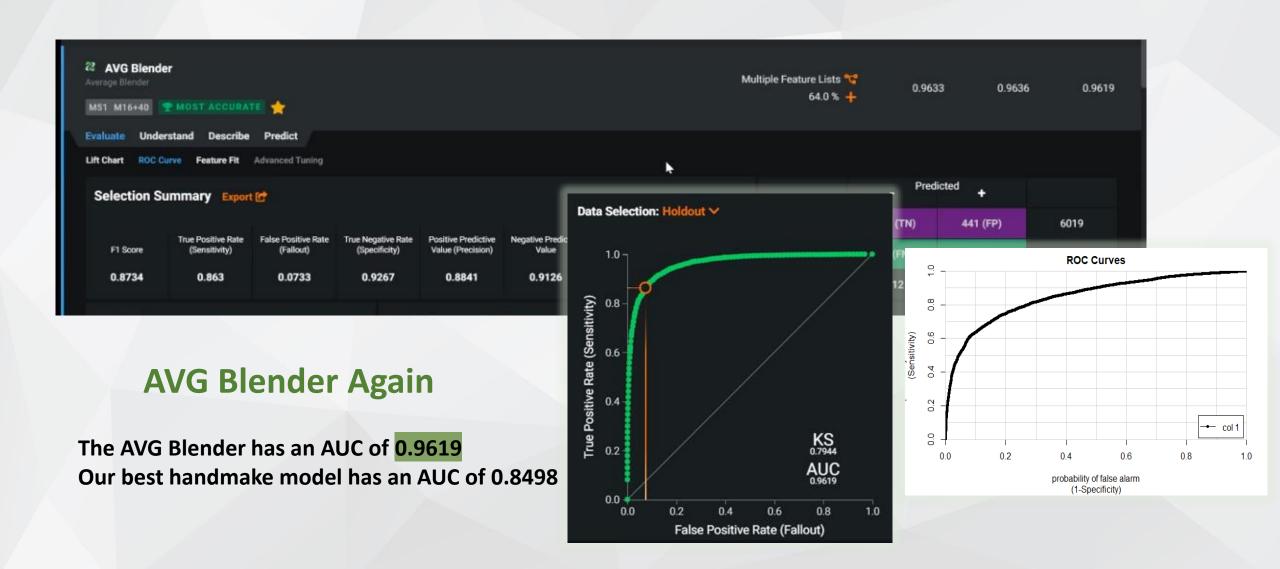
**Creating Polynomial Variables** 

**VIF Test** 

The AVG Blender has 15 variables. Our best model has 14 variables.

# 3

### **DataRobot-Comparison of Classification Models**

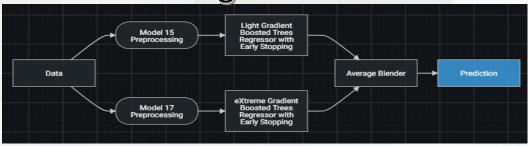


# 4

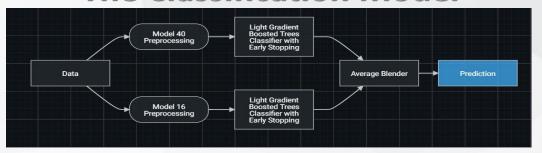
### **DataRobot-Overfitting issue**

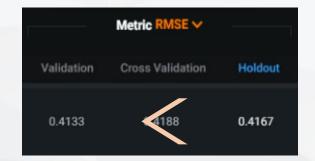
### The AVG Blender

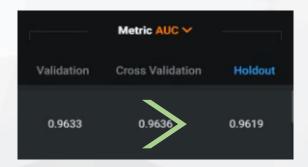
The Regression Model



#### The Classification Model

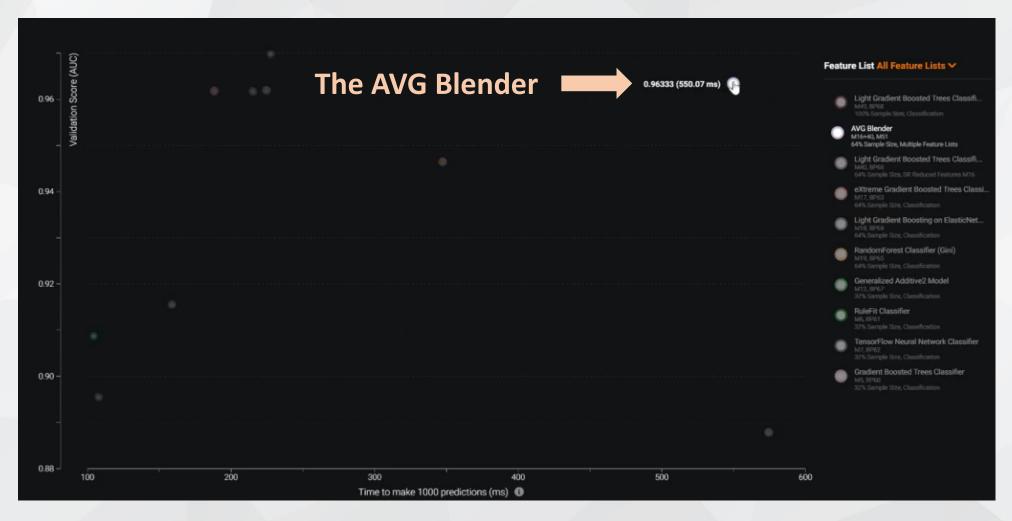






**Help Prevent Overfitting Issue** 

## 5 DataRobot-Time Consumed



# Time Consuming



## **Conclusion & Reflection**













**THANKS** 

**For Watching**