



# WINE QUALITY PREDICTION MODEL

BY: DANNY BOUDAGIAN, OLIVIA DESSO, AJDIN KUCEVIC, LUCAS  
MCDONNELL

# CONTENT

- IS THE MODEL APPROPRIATELY MOTIVATED BY THE BUSINESS USE CASE?
- ARE FEATURES TRANSFORMED APPROPRIATELY FOR MODELING?
- DOES THE PRESENTATION SHOW ENOUGH SUMMARY STATISTICS, BACKGROUND INFORMATION, AND FIGURES TO DESCRIBE THE DATA?
- IS THE FIRST MODEL APPROPRIATELY ESTIMATED?
- IS THE SECOND MODEL APPROPRIATELY ESTIMATED?
- DOES THE STUDENT PROPERLY EVALUATE THE PERFORMANCE OF THE MODEL, APPROPRIATELY FOR THE QUESTION AND DATA?
- DOES THE STUDENT COMMENT ON THE MODEL FIT AND MAKE A RECOMMENDATION REGARDING THE MOTIVATION FOR THE BUSINESS USE CASE OF THE ANALYTICS MODEL?
- ARE THE OVERALL AESTHETICS OF THE SLIDES PRESENTABLE TO A BUSINESS ENVIRONMENT?

# MOTIVATION/BUSINESS VALUE AND VARIABLES

Which variables are most significant in predicting the wine quality?

• MSE

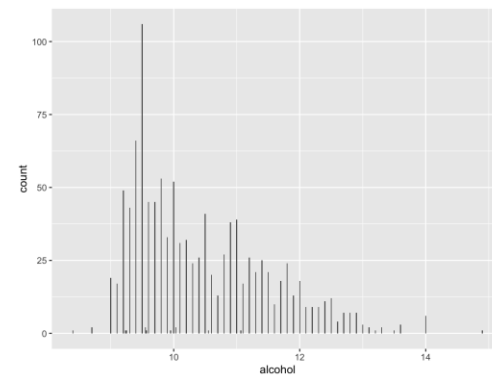
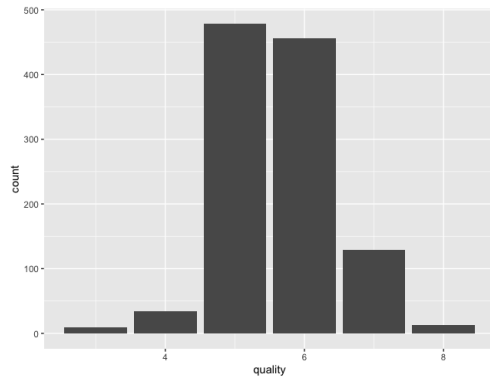
- WE WANT TO DETERMINE WHAT KINDS OF FACTORS GO INTO WINE QUALITY SO THAT PRODUCERS WILL BE ABLE TO EVALUATE HOW TO MAKE A HIGHER QUALITY WINE
- WINE INDUSTRY CURRENTLY VALUED AT **\$340 BILLION**
- THE INDUSTRY PRESERVES AGRICULTURAL LAND, AMERICAN JOBS, ATTRACTS TOURISM, AND GENERATES TAXES
- BEING ABLE TO MAKE A HIGH-QUALITY WINE CAN SIGNIFICANTLY INCREASE PROFITABILITY
- WE ARE TRYING TO UNDERSTAND AND PREDICT THE IMPACT THAT VARIABLES HAVE ON THE QUALITY OF WINE
- WE ARE USING PREDICTORS SUCH AS PH LEVEL, ALCOHOL LEVEL, DENSITY, AND CITRIC ACID LEVELS, ETC

# RAW DATA

```
> colSums(is.na(wine))
fixed.acidity      volatile.acidity      citric.acid      residual.sugar
           0              0              0              0
chlorides free.sulfur.dioxide total.sulfur.dioxide      density
           0              0              0              0
pH          sulphates          alcohol          quality
           0              0              0              0
```

```
## Rows: 1,599
## Columns: 12
## $ fixed.acidity      <dbl> 7.4, 7.8, 7.8, 11.2, 7.4, 7.4, 7.9, 7.3, 7.8, 7.5...
## $ volatile.acidity  <dbl> 0.700, 0.880, 0.760, 0.280, 0.700, 0.660, 0.600, ...
## $ citric.acid       <dbl> 0.00, 0.00, 0.04, 0.56, 0.00, 0.00, 0.06, 0.00, 0...
## $ residual.sugar    <dbl> 1.9, 2.6, 2.3, 1.9, 1.9, 1.8, 1.6, 1.2, 2.0, 6.1,...
## $ chlorides         <dbl> 0.076, 0.098, 0.092, 0.075, 0.076, 0.075, 0.069, ...
## $ free.sulfur.dioxide <dbl> 11, 25, 15, 17, 11, 13, 15, 15, 9, 17, 15, 17, 16...
## $ total.sulfur.dioxide <dbl> 34, 67, 54, 60, 34, 40, 59, 21, 18, 102, 65, 102,...
## $ density           <dbl> 0.9978, 0.9968, 0.9970, 0.9980, 0.9978, 0.9978, 0...
## $ pH                <dbl> 3.51, 3.20, 3.26, 3.16, 3.51, 3.51, 3.30, 3.39, 3...
## $ sulphates         <dbl> 0.56, 0.68, 0.65, 0.58, 0.56, 0.56, 0.46, 0.47, 0...
## $ alcohol           <dbl> 9.4, 9.8, 9.8, 9.8, 9.4, 9.4, 9.4, 10.0, 9.5, 10...
## $ quality           <int> 5, 5, 5, 6, 5, 5, 5, 7, 7, 5, 5, 5, 5, 5, 5, 7...
```

- 12 COLUMNS (11 CONTINUOUS VARIABLES AND 1 CATEGORICAL ONE)
- 1,599 ROWS
- NO DATA CLEANING REQUIRED SINCE THERE WERE NO MISSING VARIABLES IN THE DATASET
- WE USED 70-30 TRAINING-TEST SPLIT FOR OUR DATASET ON ALL MODELS
- OVERALL, WE BELIEVED THAT THESE VARIABLES SEEMED TO BE STRONG INDICATORS FOR PREDICTING THE QUALITY OF WINE



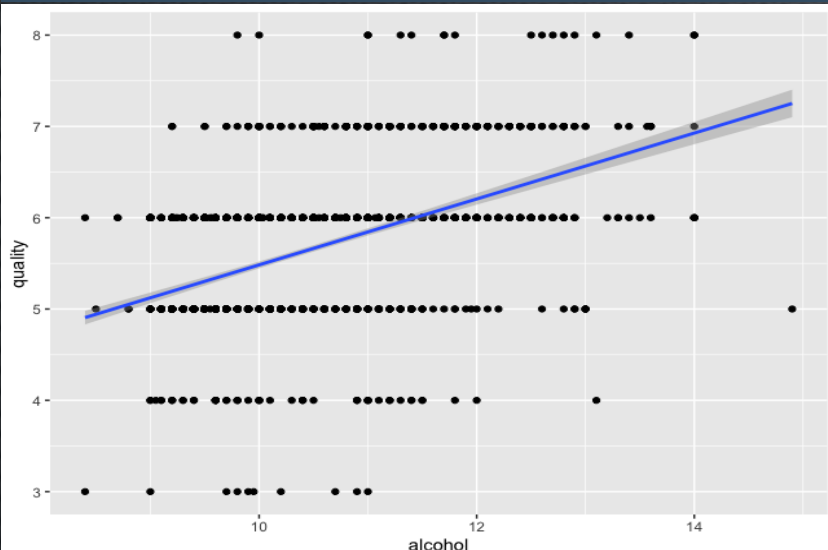
# SUMMARY STATISTICS

- THE AVERAGE WINE QUALITY CAN BE SEEN TO BE CLOSER TO 6, WITH ANYTHING ABOVE IT MEANING GOOD WINE, WHILE ANYTHING BELOW THE NUMBER 6 IS CONSIDERED BAD WINE

```
## fixed.acidity  volatile.acidity  citric.acid  residual.sugar
## Min.   : 4.60  Min.   :0.1200  Min.   :0.000  Min.   : 0.900
## 1st Qu.: 7.10  1st Qu.:0.3900  1st Qu.:0.090  1st Qu.: 1.900
## Median : 7.90  Median :0.5200  Median :0.260  Median : 2.200
## Mean   : 8.32  Mean   :0.5278  Mean   :0.271  Mean   : 2.539
## 3rd Qu.: 9.20  3rd Qu.:0.6400  3rd Qu.:0.420  3rd Qu.: 2.600
## Max.   :15.90  Max.   :1.5800  Max.   :1.000  Max.   :15.500
## chlorides      free.sulfur.dioxide  total.sulfur.dioxide  density
## Min.   :0.01200  Min.   : 1.00  Min.   : 6.00  Min.   :0.9901
## 1st Qu.:0.07000  1st Qu.: 7.00  1st Qu.: 22.00  1st Qu.:0.9956
## Median :0.07900  Median :14.00  Median : 38.00  Median :0.9968
## Mean   :0.08747  Mean   :15.87  Mean   : 46.47  Mean   :0.9967
## 3rd Qu.:0.09000  3rd Qu.:21.00  3rd Qu.: 62.00  3rd Qu.:0.9978
## Max.   :0.61100  Max.   :72.00  Max.   :289.00  Max.   :1.0037
## pH             sulphates      alcohol      quality
## Min.   :2.740  Min.   :0.3300  Min.   : 8.40  Min.   :3.000
## 1st Qu.:3.210  1st Qu.:0.5500  1st Qu.: 9.50  1st Qu.:5.000
## Median :3.310  Median :0.6200  Median :10.20  Median :6.000
## Mean   :3.311  Mean   :0.6581  Mean   :10.42  Mean   :5.636
## 3rd Qu.:3.400  3rd Qu.:0.7300  3rd Qu.:11.10  3rd Qu.:6.000
## Max.   :4.010  Max.   :2.0000  Max.   :14.90  Max.   :8.000
```

# LINEAR REGRESSION (BASELINE)

```
Residual standard error: 0.648 on 1587 degrees of freedom  
Multiple R-squared: 0.3606, Adjusted R-squared: 0.3561  
F-statistic: 81.35 on 11 and 1587 DF, p-value: < 2.2e-16
```



```
Estimate Std. Error t value Pr(>|t|)  
(Intercept) 2.197e+01 2.119e+01 1.036 0.3002  
fixed.acidity 2.499e-02 2.595e-02 0.963 0.3357  
volatile.acidity -1.084e+00 1.211e-01 -8.948 < 2e-16 ***  
citric.acid -1.826e-01 1.472e-01 -1.240 0.2150  
residual.sugar 1.633e-02 1.500e-02 1.089 0.2765  
chlorides -1.874e+00 4.193e-01 -4.470 8.37e-06 ***  
free.sulfur.dioxide 4.361e-03 2.171e-03 2.009 0.0447 *  
total.sulfur.dioxide -3.265e-03 7.287e-04 -4.480 8.00e-06 ***  
density -1.788e+01 2.163e+01 -0.827 0.4086  
pH -4.137e-01 1.916e-01 -2.159 0.0310 *  
sulphates 9.163e-01 1.143e-01 8.014 2.13e-15 ***  
alcohol 2.762e-01 2.648e-02 10.429 < 2e-16 ***  
---  
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

- WE USE LINEAR REGRESSION AS OUR BASELINE MODEL SINCE IT'S EASY TO COMPREHEND AS WELL AS IT BEING COMPUTATIONALLY INEXPENSIVE
- WE FIND THAT THE MOST STATISTICALLY SIGNIFICANT VARIABLES SEEM TO BE BOTH ALCOHOL CONTENT AND THE LEVEL OF SULFATE DIOXIDE GAS IN THE WINE BOTTLE THAT SEEM TO HAVE THE BIGGEST IMPACT ON THE QUALITY OF WINE AT A 100% CONFIDENCE LEVEL
- THE ADJUSTED  $R^2$  IS QUITE LOW (0.36) MEANING THAT OUR PREDICTORS OVERALL DON'T DO A GREAT JOB IN EXPLAINING THE VARIANCE FOUND IN THE QUALITY OF WINE AS WELL NOT BEING ABLE TO PREDICT THE MODEL AS WELL AS WE THOUGHT IT WOULD
- NOT THE START WE WANTED BUT WE CAN DEFINITELY IMPROVE THIS MODEL WITH OUR NEXT TWO MODELS IN BOTH ACCURACY AND PREDICTION IN ORDER TO HELP WINE MAKERS IMPROVE THE QUALITY OF THEIR PRODUCT TO BOTH REGULAR CUSTOMERS AND TO WINE CONNOISSEURS

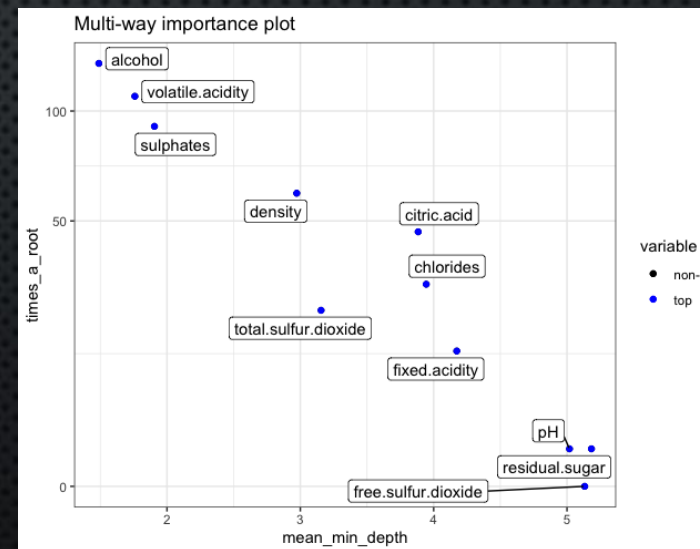
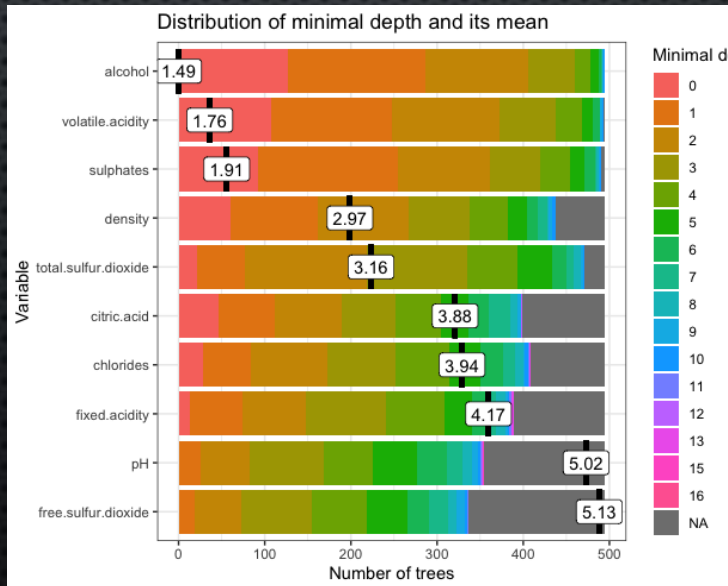
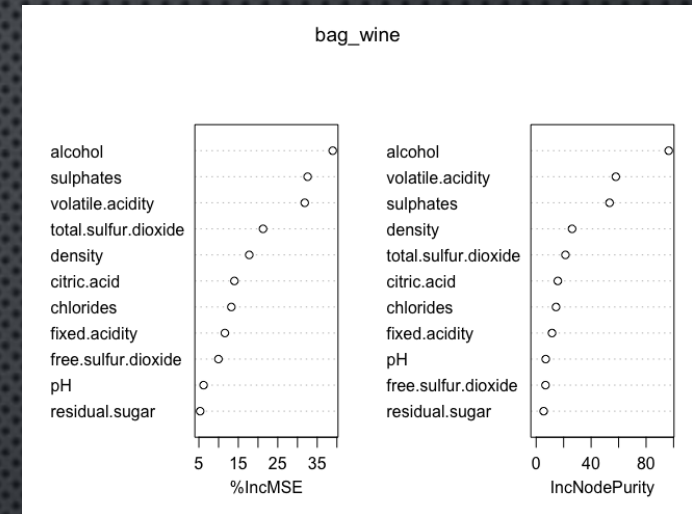
# RANDOM FOREST

- WE USE AN ENSEMBLE METHOD HERE IN ORDER TO STOP THE OVERFITTING THAT IS OCCURRING IN OUR MODEL SINCE OVERFITTING ALSO LEADS TO HIGH VARIANCE
- BOTH PLOTS (MINIMAL DEPTH DISTRIBUTION AND MULTI-WAY IMPORTANCE PLOT) HELP US DETERMINE THE IMPORTANCE OF EACH VARIABLE
- AS SEEN WITH THE LINEAR REGRESSION MODEL, WE SEE THAT OUR MODELS THAT ARE MOST SIGNIFICANT ARE ALCOHOL, VOLATILE ACIDITY, AND SULPHATES WHERE THIS CAN BE CONCLUDED SINCE THEY REQUIRE THE LEAST AMOUNT OF TREE DEPTH WHICH MEANS IN A REGRESSION TREE, THEY ARE SEEN TO BE AT THE VERY TOP OF THE TREE

```

%IncMSE  IncNodePurity
fixed.acidity  11.592998  11.550449
volatile.acidity  31.835956  57.982307
citric.acid  14.011371  15.770396
residual.sugar  5.323692  5.462965
chlorides  13.239933  14.464011
free.sulfur.dioxide  9.963321  6.805585
total.sulfur.dioxide  21.267020  21.318083
density  17.772679  26.138066
pH  6.197876  6.971193
sulphates  32.571684  53.400127
alcohol  38.934625  96.449858
> # importance plot
> varImpPlot(bag_wine)
>

```



# RANDOM FOREST

```
0 # using random forest function with mtry = p
1 bag_wine <- randomForest(quality ~ .,
2                       data = wine_train,
3                       ntree = 500,
4                       mtry = 3,
5                       nodesize = 120,
6                       err.rate = 0.1,
7                       importance = TRUE)
```

## Before

```
Call:
randomForest(formula = quality ~ ., data = wine_train, ntree = 500, mtry = 3, importance = TRUE)

Type of random forest: regression
Number of trees: 500
No. of variables tried at each split: 3

Mean of squared residuals: 0.3346744
% Var explained: 47.98
> RMSE(predicted_wine_train, wine_train$quality)
[1] 0.262253
> RMSE(predicted_wine_test, wine_test$quality) # RMSE = 2.11, 2.11k in this case
[1] 0.6049808
> |
```

## After

```
Call:
randomForest(formula = quality ~ ., data = wine_train, ntree = 500, mtry = 3, nodesize = 120, err.rate = 0.1, importance = TRUE)

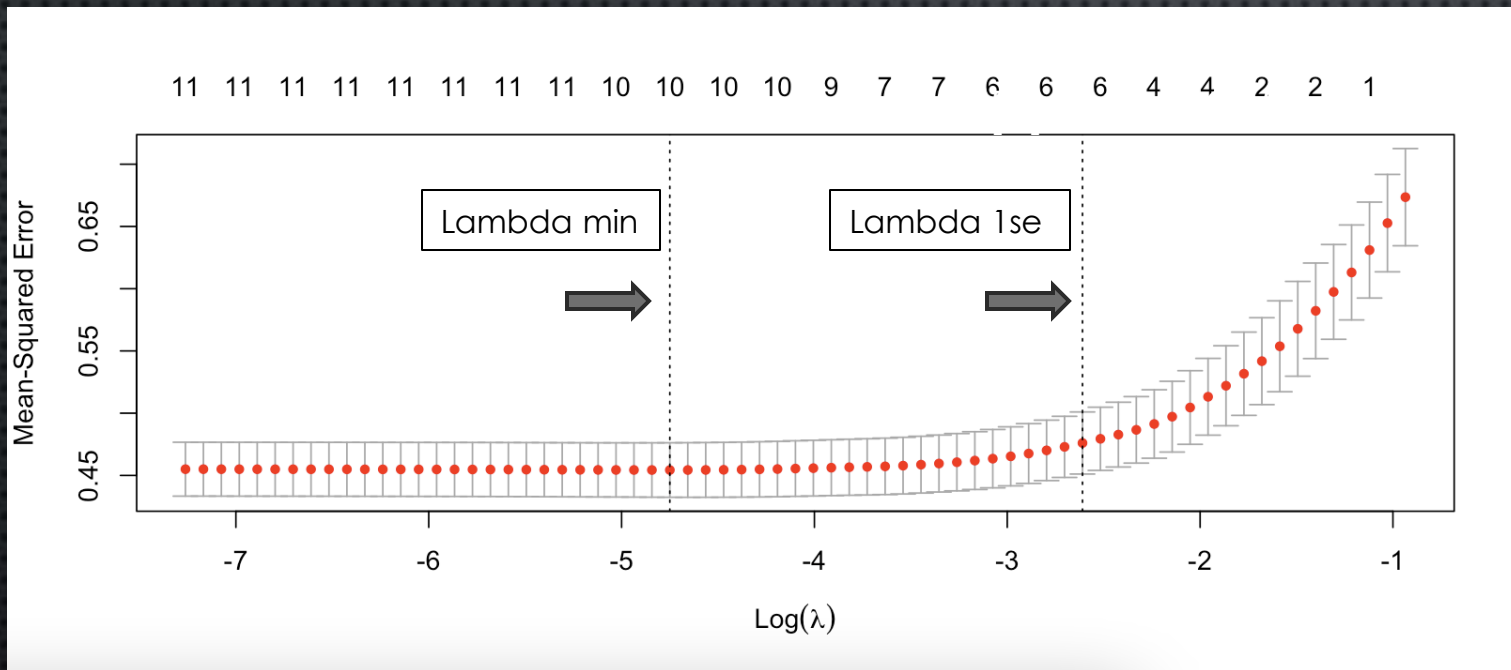
Type of random forest: regression
Number of trees: 500
No. of variables tried at each split: 3

Mean of squared residuals: 0.3956377
% Var explained: 38.5
> plot(bag_wine)
> RMSE(predicted_wine_train, wine_train$quality)
[1] 0.5704195
> RMSE(predicted_wine_test, wine_test$quality) # RMSE = 2.11, 2.11k in this case
[1] 0.6514995
> |
```

- WE USE PARAMETERS SUCH AS NODESIZE AND ERR.RATE TO REDUCE THE VARIANCE AND OVERFITTING OF OUR DATASET ALONG WITH THE USUAL IMPORTANCE, NTREE, AND MTRY PARAMETERS THAT ARE OFTEN USED IN BAGGING.
- THE IMPORTANCE OF NODESIZE IS FAIRLY UNDERESTIMATED AS IT HELPED REDUCE THE OVERFITTING MASSIVELY
- NODESIZE- DETERMINES THE MINIMUM NUMBER OF OBSERVATIONS IN EACH TERMINAL LEAF NODE
- THE HIGHER THE NODESIZE THE FEWER LEAF NODES WE HAVE, WHICH REDUCES THE COMPLEXITY OF THE MODEL
- WHEN PLOTTING THE RANDOM FOREST FUNCTION, WE SEE THAT AROUND 100 TREES GAVE USE THE LEAST AMOUNT OF ERROR WITH THE ERROR BEING CONSTANT THEREAFTER
- ERR.RATE - SLIGHT DIFFERENCE, GIVES US A SMALLER DIFFERENCE IN ERROR SIZE BUT NOTHING DRASTIC



# LASSO REGRESSION



```
# Estimating the model
```

```
lasso_mod = cv.glmnet(quality ~ .,  
                      data = wine_train,  
                      alpha = 1)
```

```
# Lambda
```

```
print(lasso_mod$lambda.min) [1] 0.005272493  
print(lasso_mod$lambda.1se) [1] 0.0859285
```

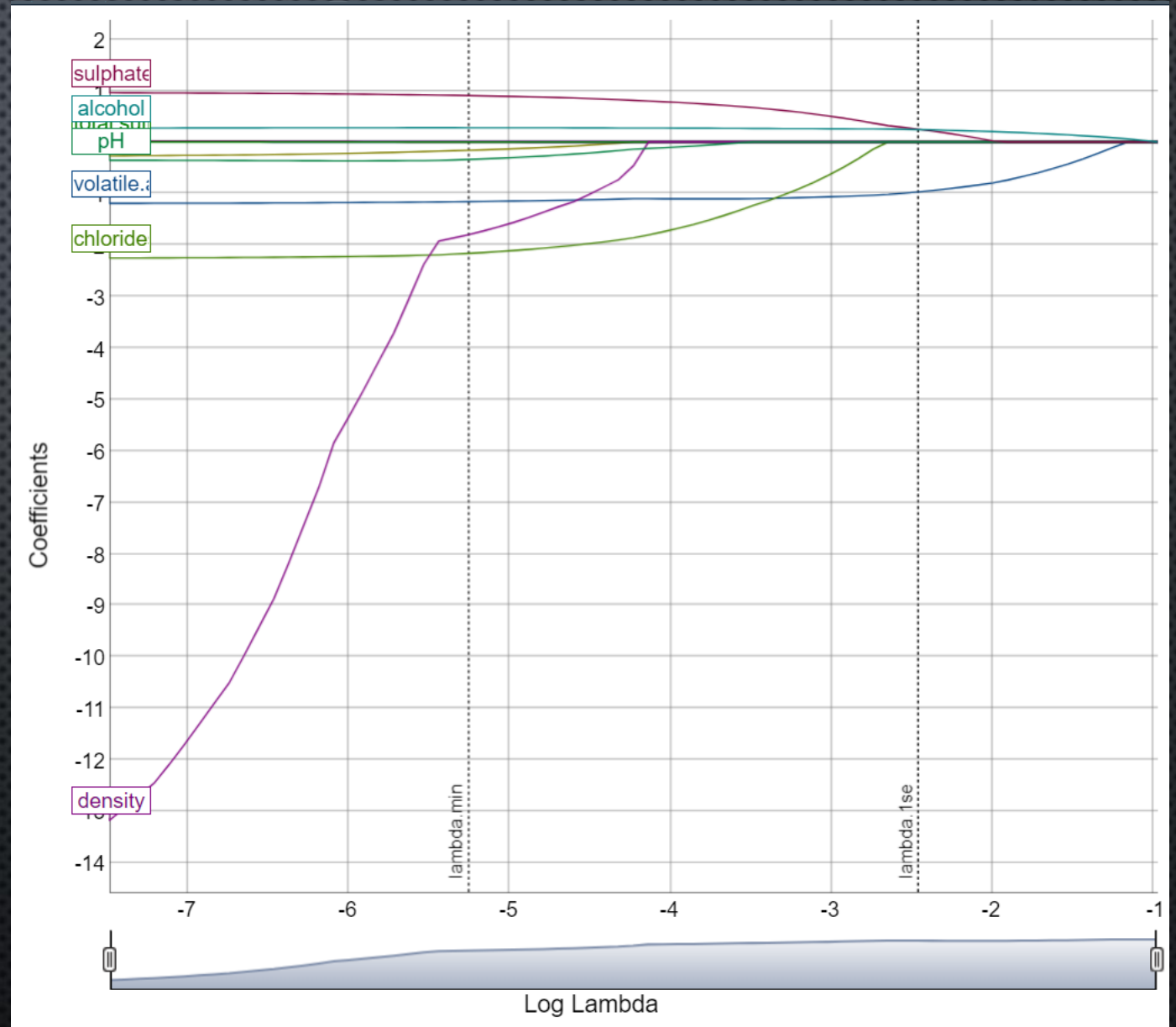
```
# Coefficients
```

```
coef(lasso_mod, s = lasso_mod$lambda.1se) %>%  
  as.matrix() %>%  
  as.data.frame() %>%  
  round(3)
```

```
coef(lasso_mod, s = lasso_mod$lambda.min) %>%  
  as.matrix() %>%  
  as.data.frame() %>%  
  round(3)
```

# LASSO REGRESSION

- SIMILAR TO THE LINEAR MODEL, VARIABLES LIKE ALCOHOL, SULPHATES, AND ACIDITY ARE MOST SIGNIFICANT
- MANY VARIABLES HAVE MINIMAL SIGNIFICANCE
- LAMBDA MIN MORE APPROPRIATE



## Lambda 1se

	s1
(Intercept)	3.394
fixed.acidity	0.000
volatile.acidity	-0.996
citric.acid	0.000
residual.sugar	0.000
chlorides	0.000
free.sulfur.dioxide	0.000
total.sulfur.dioxide	0.000
density	0.000
pH	0.000
sulphates	0.285
alcohol	0.247

# COEFFICIENTS

```

lasso_coefs <- data.frame(
  lasso_min = coef(lasso_mod, s = lasso_mod$lambda.min) %>%
    round(3) %>% as.matrix(),
  lasso_1se = coef(lasso_mod, s = lasso_mod$lambda.1se) %>%
    round(3) %>% as.matrix()
) %>% rename(lasso_min = 1, lasso_1se = 2)

print(lasso_coefs)

lasso_coefs %>%
  select(lasso_min) %>%
  filter(lasso_min != 0) %>%
  nrow()
[1] 10

lasso_coefs %>%
  select(lasso_1se) %>%
  filter(lasso_1se != 0) %>%
  nrow()
[1] 4
  
```

10 vs 4  
non-zero  
predictors

## Lambda min

	s1
(Intercept)	3.467
fixed.acidity	0.000
volatile.acidity	-1.102
citric.acid	0.000
residual.sugar	0.000
chlorides	-1.742
free.sulfur.dioxide	0.000
total.sulfur.dioxide	-0.002
density	0.000
pH	-0.110
sulphates	0.783
alcohol	0.271

- lasso\_1se coefficients have more penalization, are "shrunk" closer to zero
- Therefore, more zero coefficients in the lasso\_1se model compared to lasso\_min

# PREDICTION & EVALUATION

```
#Prediction
predict_train <- predict(lasso_mod, s=lasso_mod$lambda.min,wine_train)
predict_test <- predict(lasso_mod,s=lass0_mod$lambda.min_wine_test)

#Evaluation
results_train <- wine_train %>% mutate(pred=predict_train)
RMSE(results_train$pred,results_train$quality) [1] 0.6389668
```

- Choosing Lambda:
  - Lasso\_min model has smallest value of lambda, therefore lowest training error-->prone to overfitting data
  - Lasso\_1se larger value of lambda, higher training error--> may be better for generalization
- In the context of this Lasso model, an RSME of 0.6389 means that on average, the model's predicted response values are off by approximately 0.6398 units from the actual response values.
  - Pretty good considering the values are relatively small

# CONCLUSION

- WE USED RANDOM FOREST AND LASSO REGRESSION MODELS COMPARED TO LINEAR REGRESSION AS OUR BASELINE TO ANALYZE THE DATA AND CONCLUDED THAT ALCOHOL, SULPHATES, AND VOLATILE ACIDITY WERE THE MOST IMPORTANT VARIABLES TO PREDICTING QUALITY OF WINE
- THE USE OF RANDOM FOREST WAS DUE TO ITS EASE OF USE, FLEXIBILITY FOR BOTH CLASSIFICATION AND REGRESSION TREES, AS WELL AS BEING MUCH MORE ACCURATE THAN BAGGING IN THE PREDICTION OF OUR MODELS
- LASSO REGRESSION ALLOWS US TO UNDERSTAND WHICH VARIABLES SHOULD BE CONSIDERED AS SIGNIFICANT TO THE QUALITY OF WINE AND INTERPRET IT EASILY FOR MANAGEMENT OR BUSINESS PURPOSES
- THIS KNOWLEDGE COULD HELP PRODUCERS CREATE HIGHER QUALITY OF WINE AND INCREASE PROFITABILITY FOR THE THEM