

# Supervised Learning:

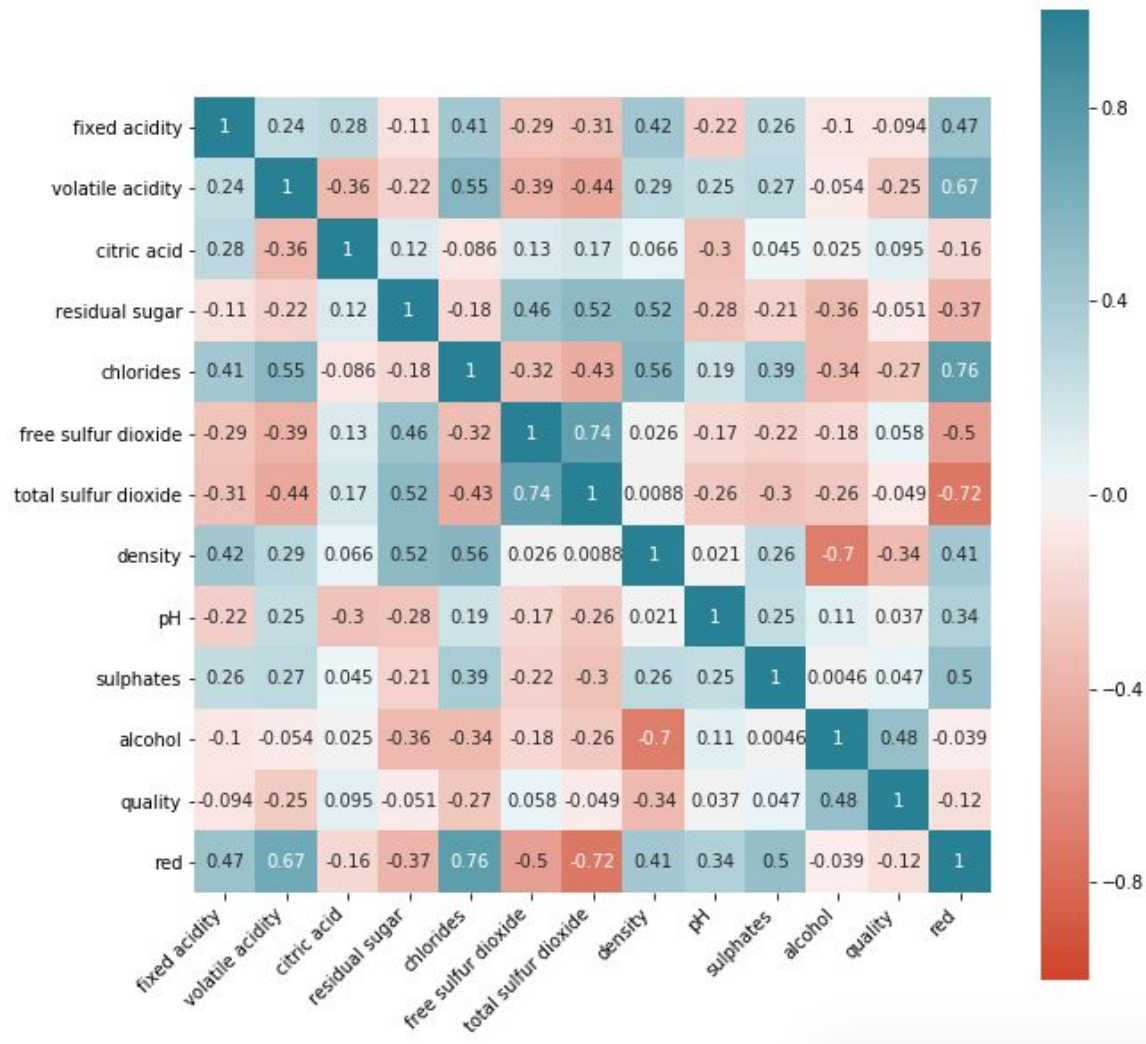
Predicting Wine  
Color

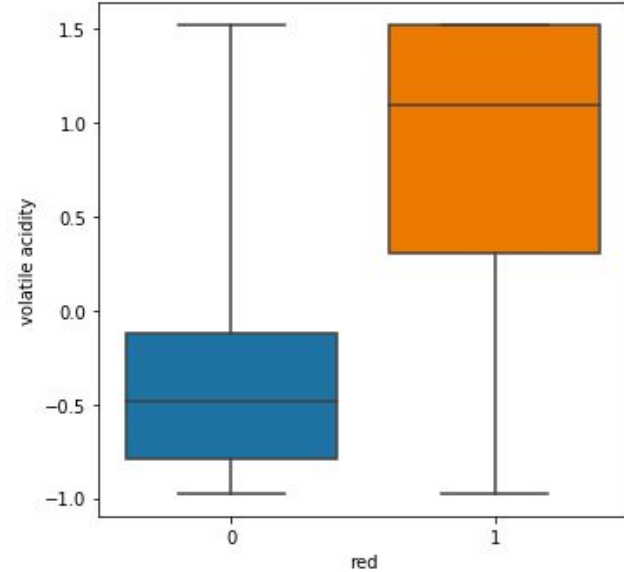
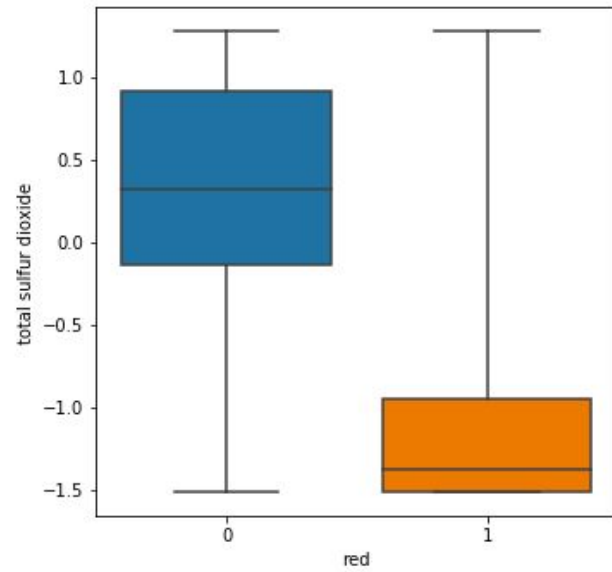
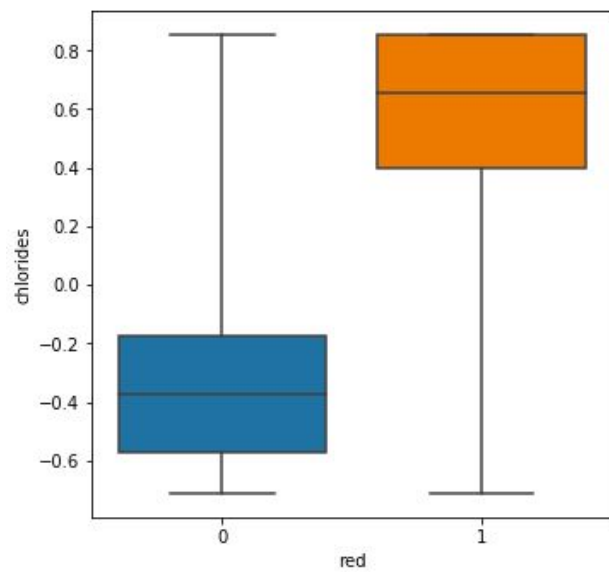
# Motivation

- Humans can differentiate wines by their color and flavor, but can we teach a computer to distinguish them by what chemicals they contain?
  - What model would be best for this classification problem?
  - What variables are most important to this classification?

# Data Source

- UC Irvine Machine Learning Repository - Wine Quality
  - Red Wine & White Wine - Imbalanced Set
  - 11 Features Regarding the Chemical Composition of the Wine
  - 1 Feature Regarding the Overall “Rating” of the Wine’s Flavor
- Modifications
  - Indicator Variable
  - Standardization
  - Winsorization





# Linear Classification

## Using 3 Variables

Score on training data: 0.987673343605547

Score on test data: 0.9779411764705882

Mean cross validation score: 0.986

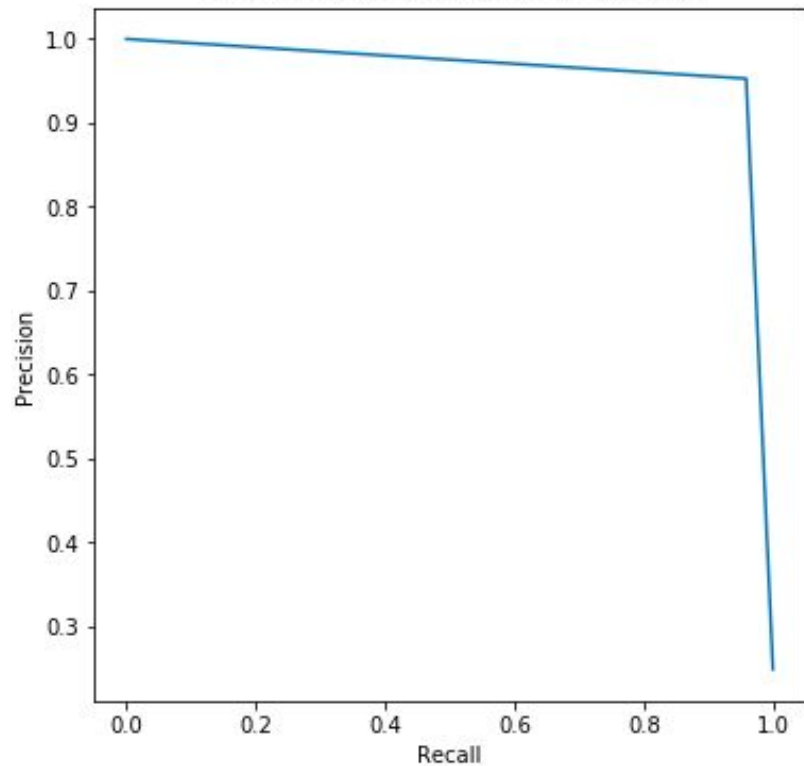
## Using 2 Variables

Score on training data: 0.9599383667180277

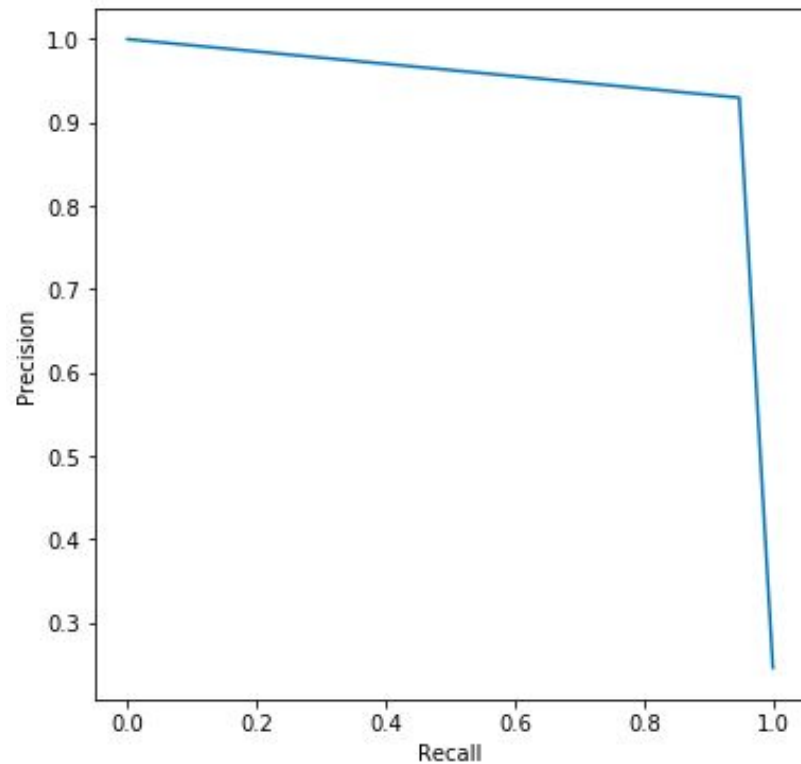
Score on test data: 0.969562243502052

Mean cross validation score: 0.958

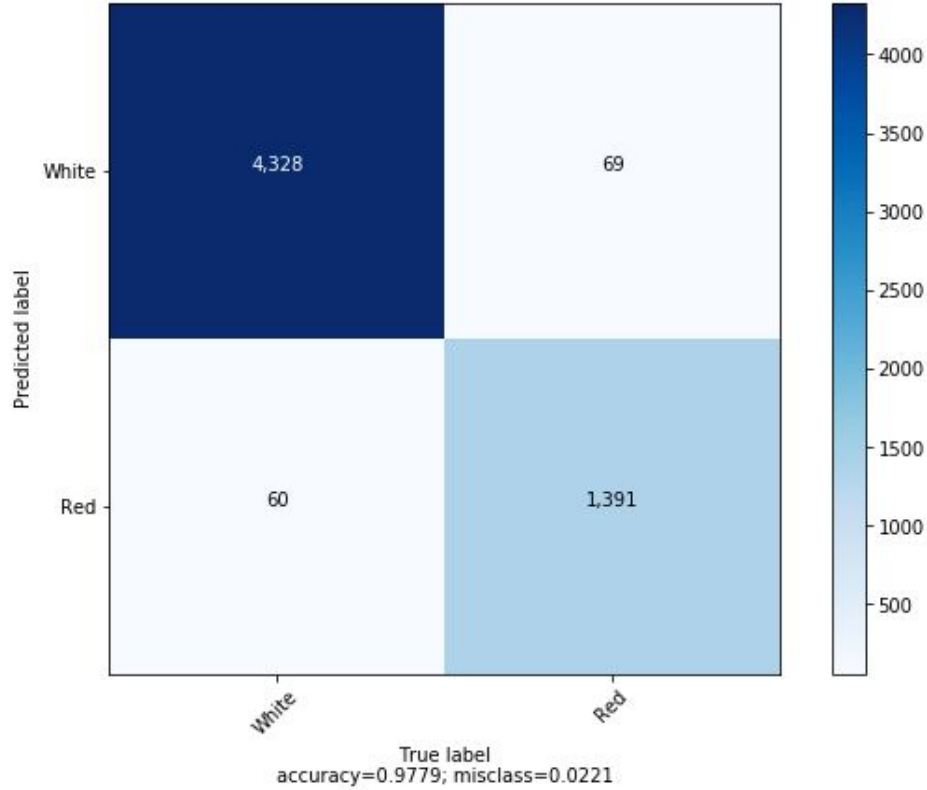
Precision-Recall Curve With 3 Variables



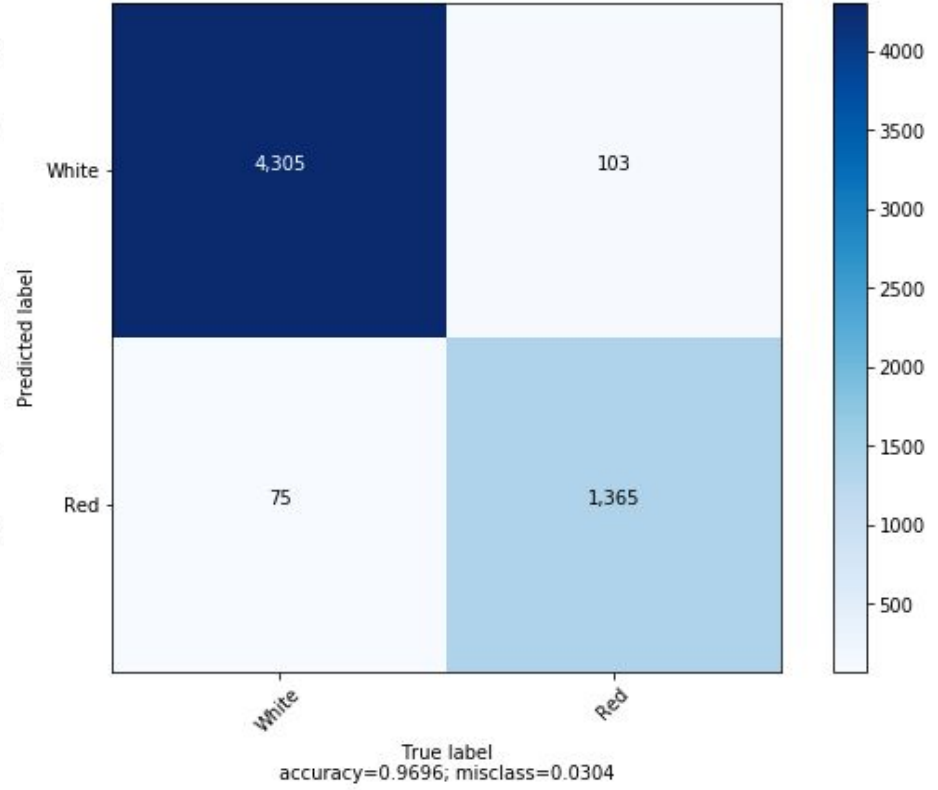
Precision-Recall Curve With 2 Variables



Confusion Matrix With 3 Variables



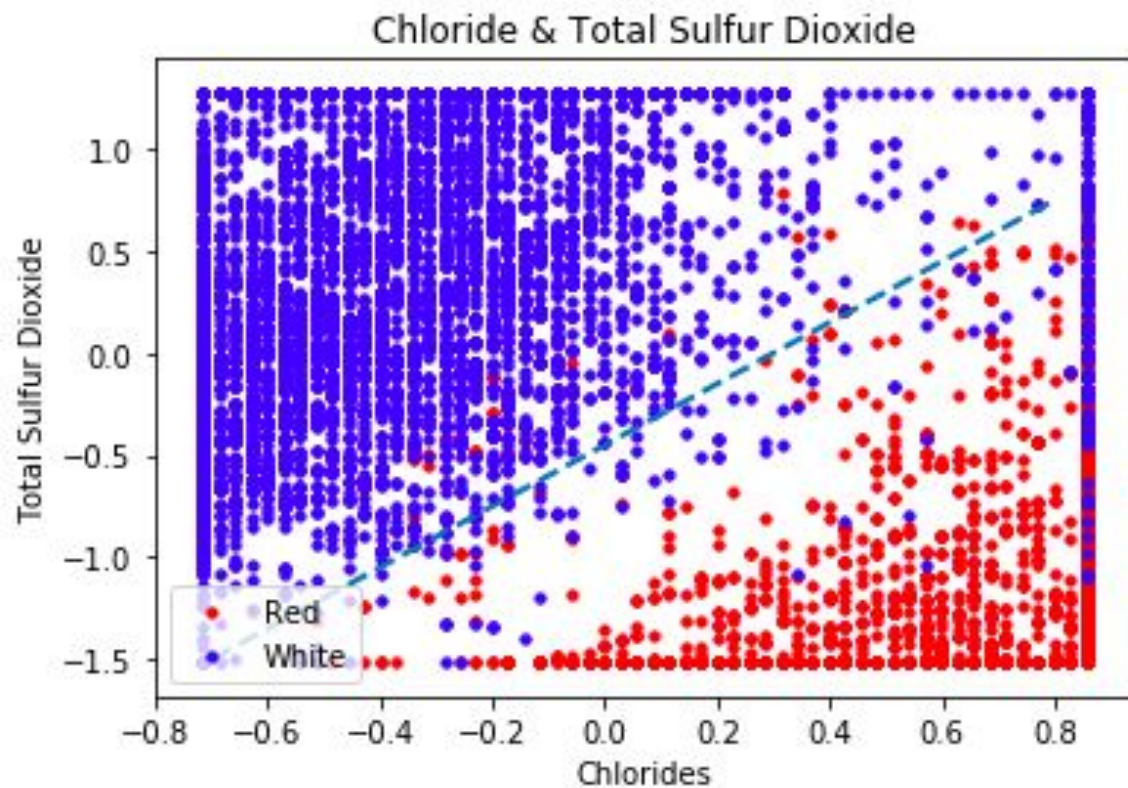
Confusion Matrix With 2 Variables





It appears that the model is overfit when I use 3 variables, so I will continue the analysis by using 2 variables – chlorides and total sulfur dioxide.

# SVC



There doesn't appear to be a margin that perfectly splits the data, so SVC might not be the best approach. What we do notice, however, is the similarity between points of the same class.

# KNN Neighbors

No Weights

Distance

K=5

Score on training data: 0.9768875192604006  
Score on test data: 0.9805061559507524  
Mean cross validation score: 0.9687553739562977

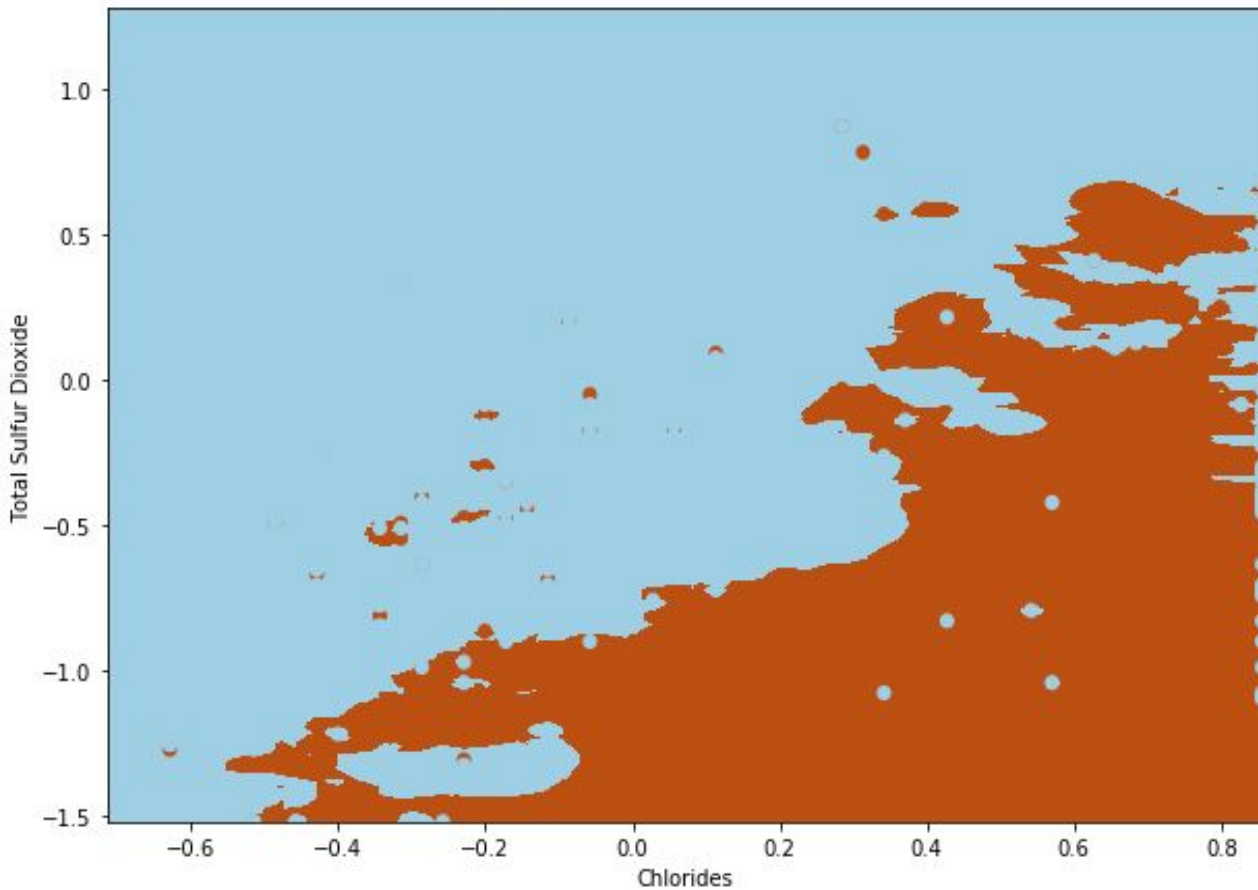
Score on training data: 0.9892141756548536  
Score on test data: 0.9904240766073872  
Mean cross validation score: 0.9682940723633564

K=10

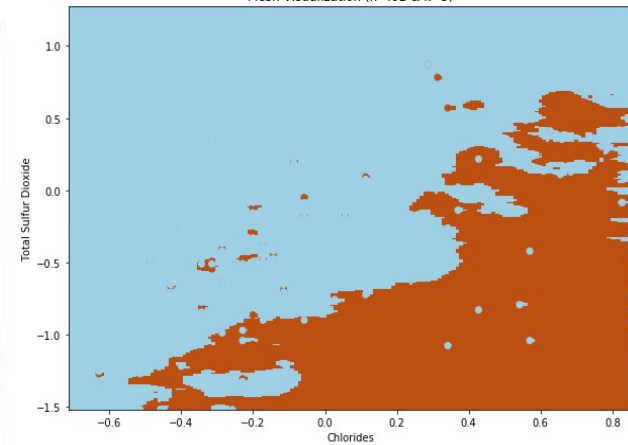
Score on training data: 0.9830508474576272  
Score on test data: 0.9752051983584131  
Mean cross validation score: 0.9687553739562977

Score on training data: 0.9953775038520801  
Score on test data: 0.9897400820793434  
Mean cross validation score: 0.9689095754130397

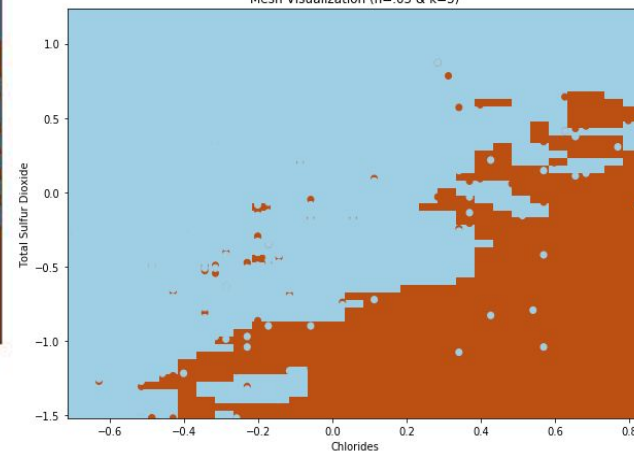
Mesh Visualization ( $h=.001$  &  $k=5$ )



Mesh Visualization ( $h=.01$  &  $k=5$ )

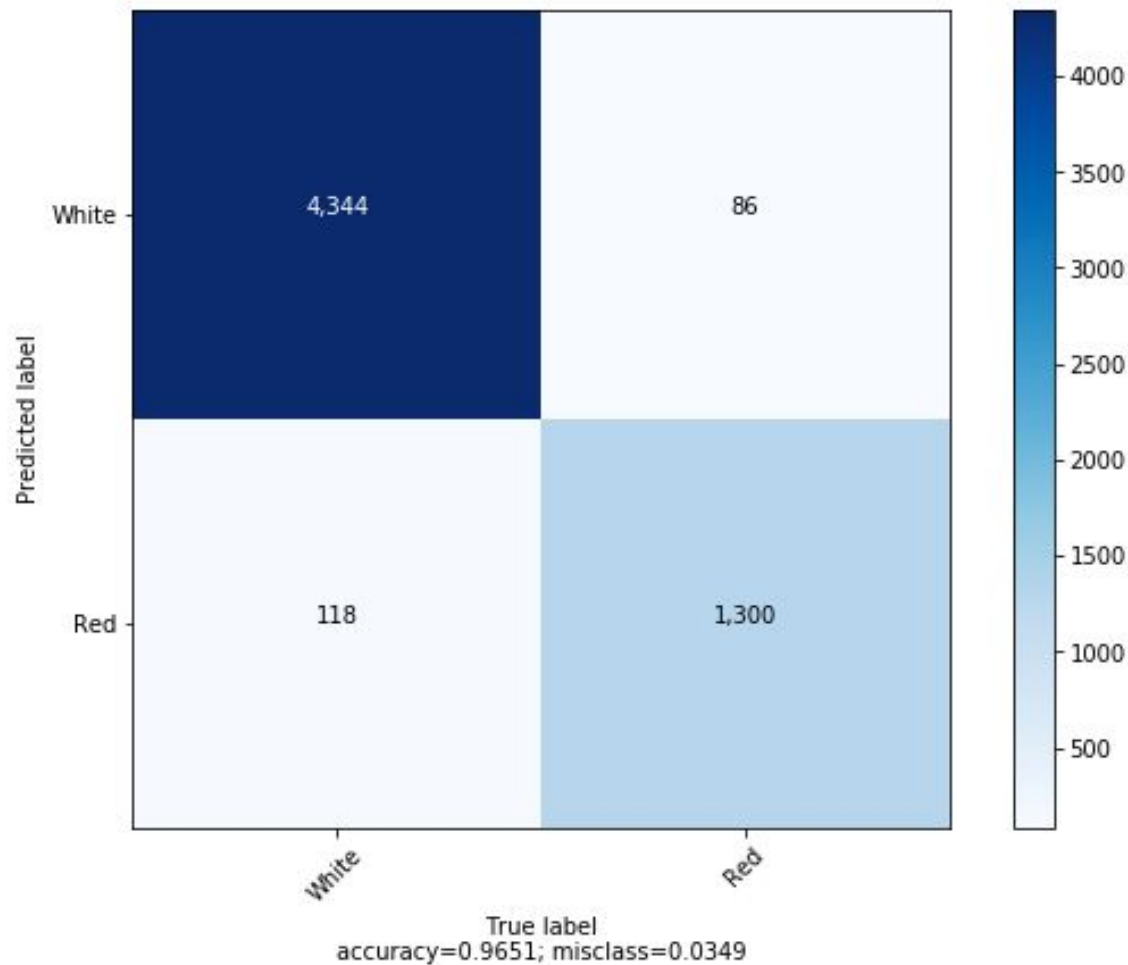


Mesh Visualization ( $h=.05$  &  $k=5$ )

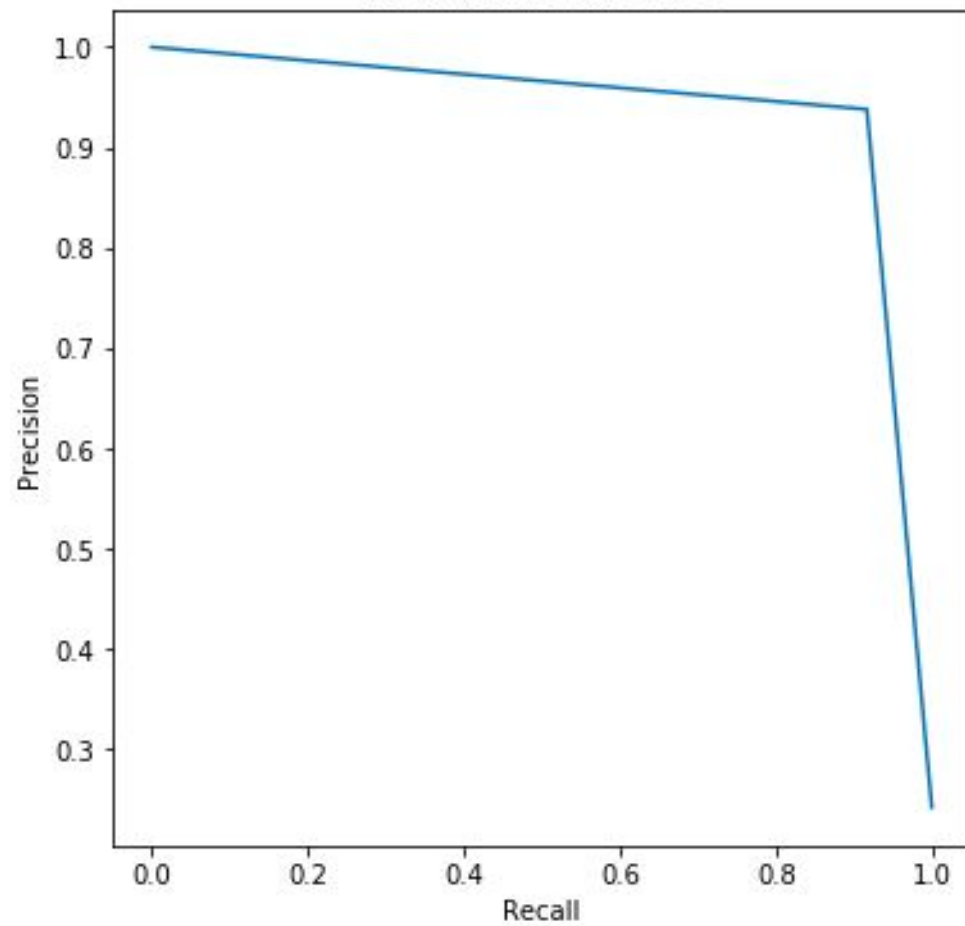


Decreasing the step size makes the mesh smoother; however, the boundaries of the graph continue to be problematic due to the standardization of the data set. In general, however, it appears red wines have less sulfur dioxide and more chloride than white wine.

Confusion Matrix

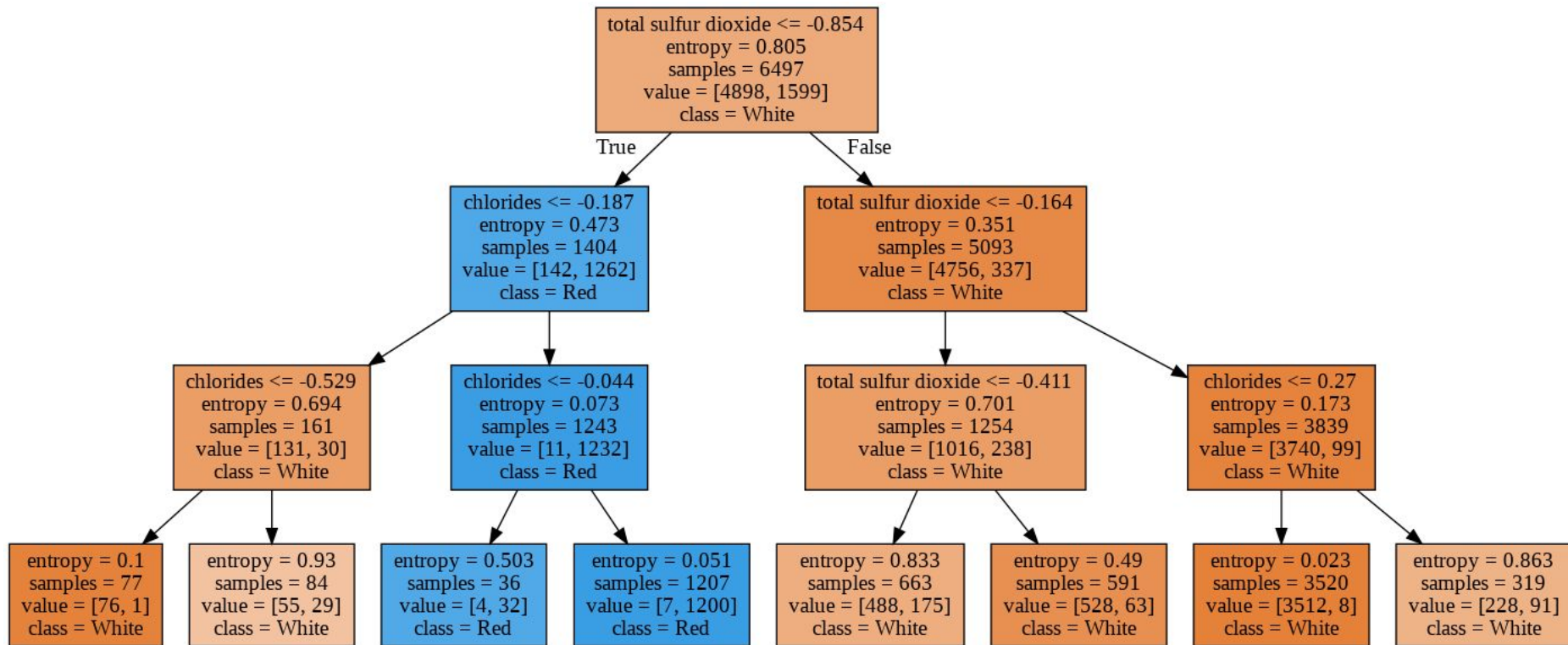


Precision-Recall Curve





# Tree Models



**Notice the criteria at each node is either related to chlorides or total sulfur dioxide. This reaffirms that these variables are preferred to the Decision Tree Model, just as they were in our previous models.**

## Decision Tree

Depth = 2

Score on training data: 0.9567901234567902  
Score on test data: 0.9562611372104325  
Mean cross validation score: 0.940278912773139

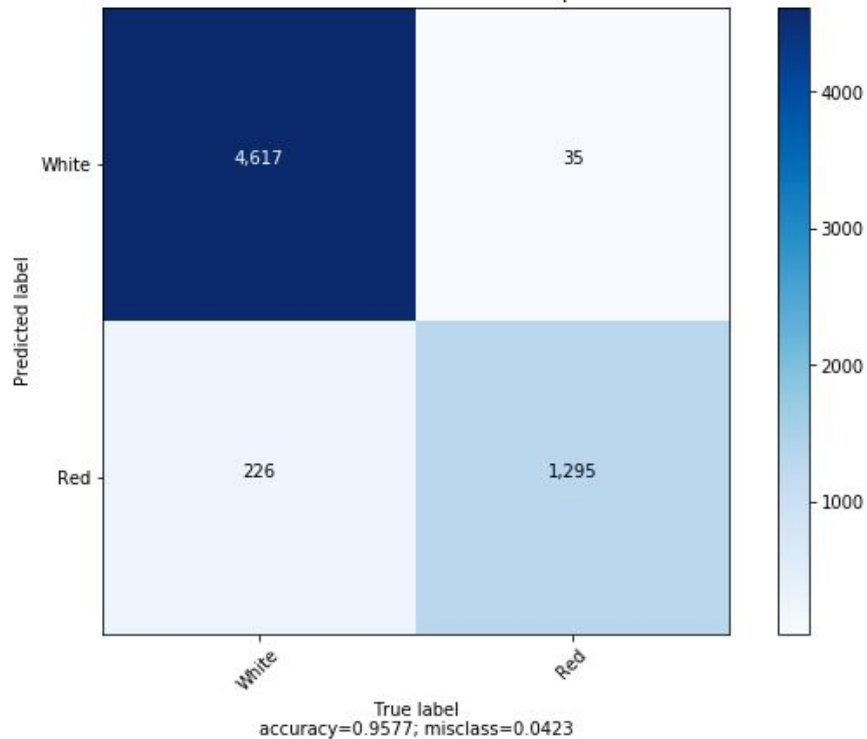
Depth = 3

Score on training data: 0.9444444444444444  
Score on test data: 0.9577190993034181  
Mean cross validation score: 0.9595202226564814

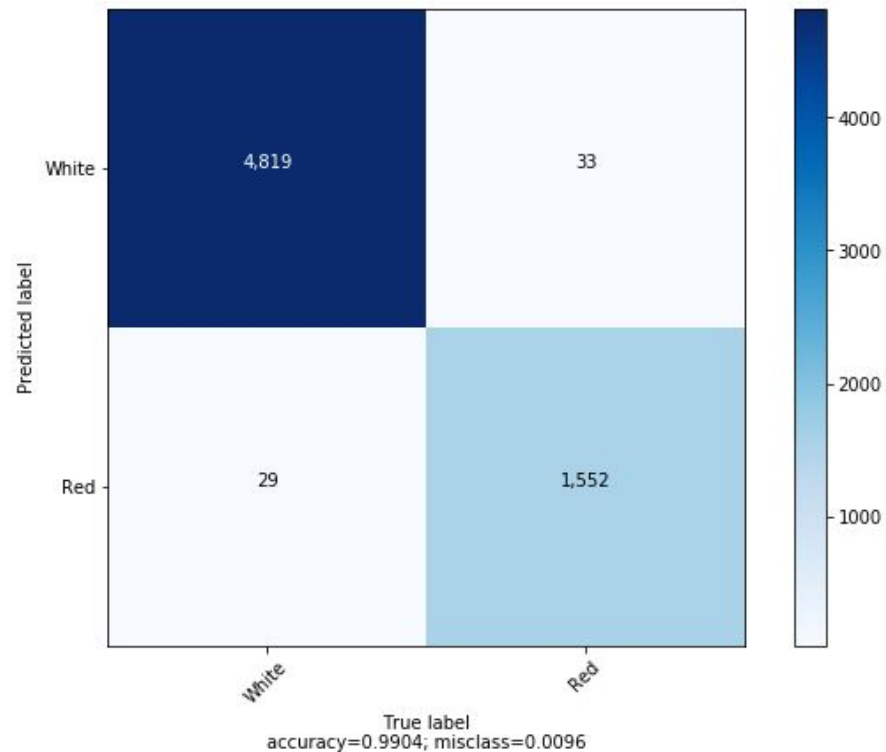
## Random Forest

Score on training data: 0.984375  
Score on test data: 0.9902067464635473  
Mean cross validation score: 0.9682934801918636

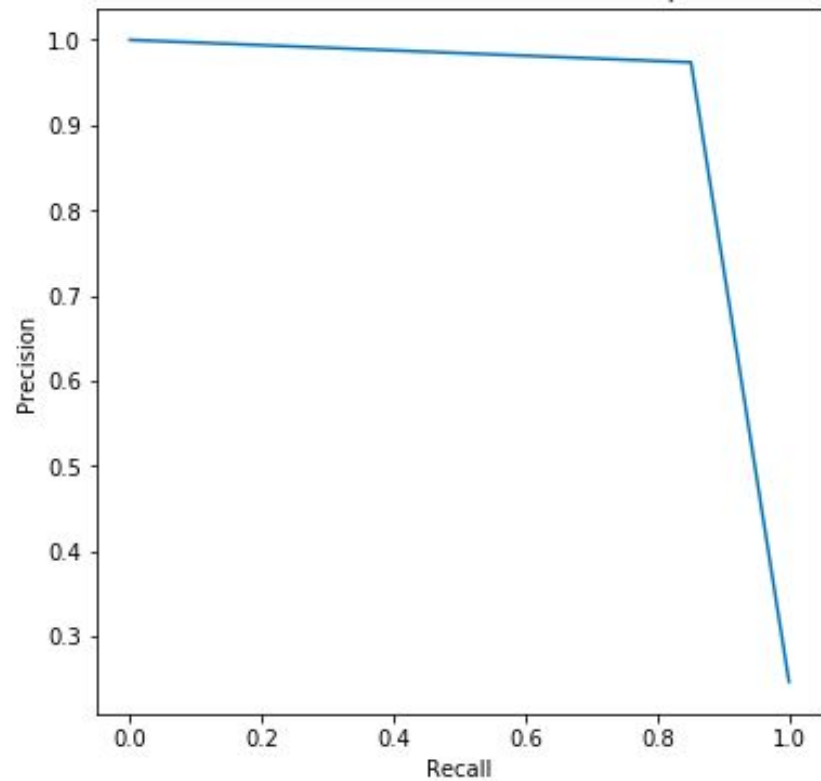
Confusion Matrix for Tree With Depth=3



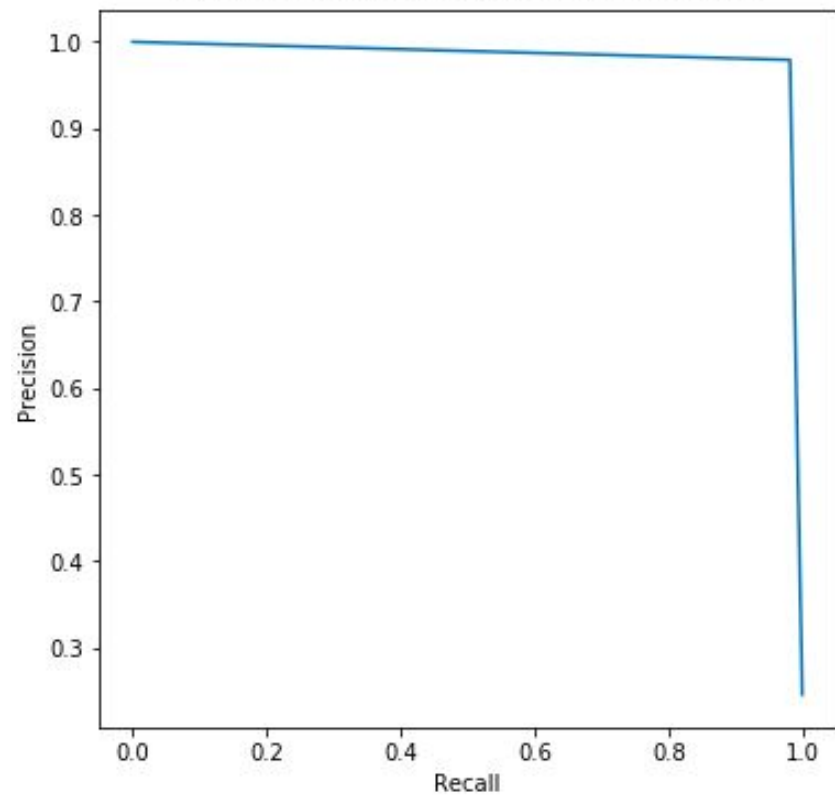
Confusion Matrix with Random Forest



Precision-Recall Curve for Tree With Depth = 3



Precision-Recall Curve with Random Forest



**The Random Forest Model seems to perfect the Decision Tree Model through repetition without costing too much time.**

# Conclusion

- At a chemical level, the differences between red wine & white wine are very clear to a machine, thus classifications are extremely effective in nearly every model
  - Chlorides and Sulfur Dioxide appear to be the most influential variables in this classification problem
- The Random Forest Model performed at the highest level of accuracy without being overfit to the data