

Wine Quality Prediction

By Taya Bhavsar



Goal of Study

Main Objective: Determine which physico-chemical properties most accurately predict red wine quality



Research Questions

Three Research Questions:

1. Can quality be predicted using only 3 chemical measurements?
2. Is alcohol content alone sufficient, or are acidity measures necessary?
3. What minimum alcohol percentage gives >50% probability of good quality?

Dataset Overview

Source: Portuguese "Vinho Verde" red wines
(UCI Wine Quality dataset)

Size: 1,599 samples, no missing values

Features: 11 continuous physicochemical properties (fixed acidity, volatile acidity, citric acid, residual sugar, chlorides, free sulfur dioxide, total sulfur dioxide, density, pH, sulphates, alcohol)

Target: Quality scores 3-8 (mean: 5.636)

Key Challenge: Imbalanced dataset (46.5% poor quality <6, 53.5% good quality ≥6)

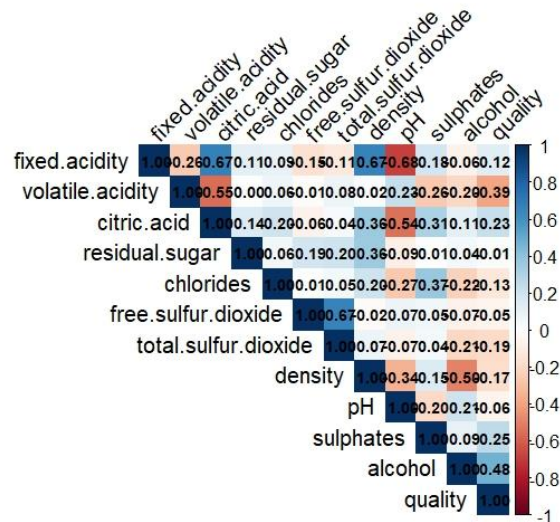


Wine Quality Correlation Matrix

Strongest Correlations (Dark Blue)

- **Fixed acidity ↔ Citric acid (0.67):**
Wines with higher fixed acidity tend to have more citric acid, reflecting natural grape composition
- **Fixed acidity ↔ Density (0.67):**
Higher acid content increases wine density
- **Fixed acidity ↔ pH (-0.68):** This negative correlation is chemically expected—more acid means lower pH (more acidic)

Wine Quality Feature Correlation Matrix



Low multicollinearity risk

Models Used

Three complementary approaches:

1. Logistic Regression

- Binary classification (Good ≥ 6 vs Bad < 6)
- Identifies most significant predictors
- Good for interpretability

2. Classification Trees

- Multi-class prediction (Low: 3-4, Medium: 5-6, High: 7-8)
- Reveals critical thresholds for alcohol and volatile acidity

3. Gradient Boosting/Random Forest

- Handles class imbalance and non-linear relationships
- Captures complex chemical interactions
- Handles complexity

Evaluation: 80/20 train/test split with 1,279 training and 320 test samples using accuracy, precision, recall, F1-score

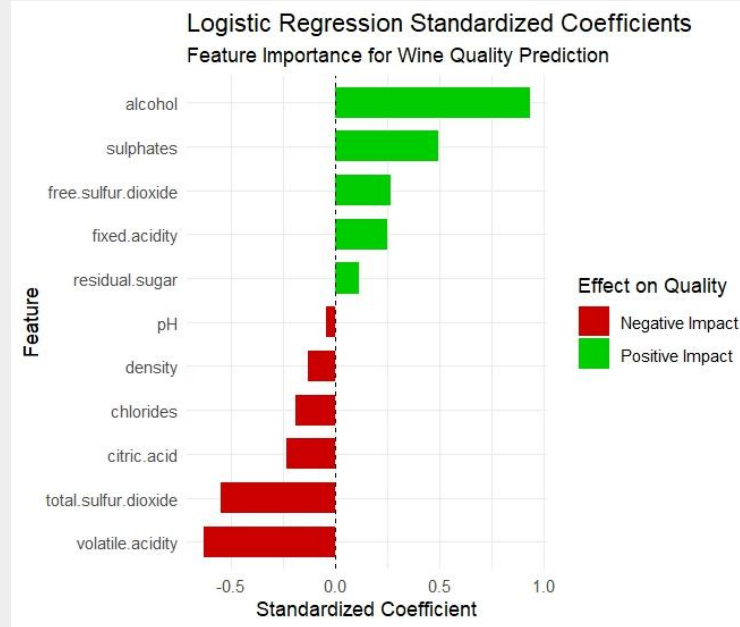
Main Findings

Training Set: 1,279 samples

Test Set: 320 samples

Top 3 Predictors Identified: Logistic Regression

1. **Alcohol content** - strongest positive predictor
2. **Volatile acidity** - strongest negative predictor (vinegar-like flavors hurt quality)
3. **Sulphates** - positive effect on quality



Binary Classification Tree

Accuracy: 74.38% - Successfully classifies 3 out of 4 wines correctly

AUC: 0.7963 - Good discriminatory power between quality classes

Balanced Performance: Sensitivity (74.25%) and Specificity (74.51%) are nearly equal, indicating the model performs well on both Good and Bad wines

```
graph TD
    Root["Good  
.48 .54  
100%"]
    Root -- "yes - alcohol < 10" --> Node1["Bad  
.66 .34  
52%"]
    Root -- "no" --> Node2["Good  
.25 .75  
48%"]
    Node1 -- "sulphates < 0.58" --> Node3["Bad  
.55 .45  
31%"]
    Node1 -- "no" --> Node4["Good  
.25 .75  
48%"]
    Node3 -- "total.sulfur.dioxide >= 51" --> Node5["Good  
.43 .57  
18%"]
    Node3 -- "no" --> Node6["Bad  
.61 .39  
7%"]
    Node5 -- "density < 1" --> Node7["Bad  
.61 .39  
7%"]
    Node5 -- "no" --> Node8["Good  
.31 .69  
11%"]
    Node7 -- "volatile.acidity >= 0.39" --> Node9["Bad  
.65 .35  
6%"]
    Node7 -- "no" --> Node10["Good  
.49 .51  
4%"]
    Node9 -- "citric.acid >= 0.01" --> Node11["Bad  
.80 .20  
22%"]
    Node9 -- "no" --> Node12["Bad  
.72 .28  
13%"]
    Node10 -- "sulphates < 0.66" --> Node13["Good  
.33 .67  
3%"]
    Node10 -- "no" --> Node14["Good  
.20 .80  
7%"]
    Node13 -- "volatile.acidity >= 0.55" --> Node15["Good  
.62 .38  
8%"]
    Node13 -- "no" --> Node16["Good  
.00 1.00  
1%"]
    Node15 -- "volatile.acidity >= 0.55" --> Node17["Bad  
1.00 .00  
1%"]
    Node15 -- "no" --> Node18["Good  
.22 .78  
22%"]
    Node17 -- "volatile.acidity >= 0.55" --> Node19["Good  
.08 .92  
17%"]
    Node17 -- "no" --> Node20["Good  
.08 .92  
17%"]
    Node2 -- "alcohol < 11" --> Node21["Good  
.34 .66  
31%"]
    Node2 -- "no" --> Node22["Good  
.24 .76  
22%"]
    Node21 -- "sulphates < 0.59" --> Node23["Bad  
.58 .42  
9%"]
    Node21 -- "no" --> Node24["Good  
.24 .76  
22%"]
    Node23 -- "volatile.acidity >= 0.35" --> Node25["Bad  
.58 .42  
9%"]
    Node23 -- "no" --> Node26["Good  
.24 .76  
22%"]
    Node25 -- "volatile.acidity >= 0.35" --> Node27["Bad  
.58 .42  
9%"]
    Node25 -- "no" --> Node28["Good  
.24 .76  
22%"]
    Node27 -- "volatile.acidity >= 0.35" --> Node29["Bad  
.58 .42  
9%"]
    Node27 -- "no" --> Node30["Good  
.24 .76  
22%"]
    Node29 -- "volatile.acidity >= 0.35" --> Node31["Bad  
.58 .42  
9%"]
    Node29 -- "no" --> Node32["Good  
.24 .76  
22%"]
    Node31 -- "volatile.acidity >= 0.35" --> Node33["Bad  
.58 .42  
9%"]
    Node31 -- "no" --> Node34["Good  
.24 .76  
22%"]
    Node33 -- "volatile.acidity >= 0.35" --> Node35["Bad  
.58 .42  
9%"]
    Node33 -- "no" --> Node36["Good  
.24 .76  
22%"]
    Node35 -- "volatile.acidity >= 0.35" --> Node37["Bad  
.58 .42  
9%"]
    Node35 -- "no" --> Node38["Good  
.24 .76  
22%"]
    Node37 -- "volatile.acidity >= 0.35" --> Node39["Bad  
.58 .42  
9%"]
    Node37 -- "no" --> Node40["Good  
.24 .76  
22%"]
    Node39 -- "volatile.acidity >= 0.35" --> Node41["Bad  
.58 .42  
9%"]
    Node39 -- "no" --> Node42["Good  
.24 .76  
22%"]
    Node41 -- "volatile.acidity >= 0.35" --> Node43["Bad  
.58 .42  
9%"]
    Node41 -- "no" --> Node44["Good  
.24 .76  
22%"]
    Node43 -- "volatile.acidity >= 0.35" --> Node45["Bad  
.58 .42  
9%"]
    Node43 -- "no" --> Node46["Good  
.24 .76  
22%"]
    Node45 -- "volatile.acidity >= 0.35" --> Node47["Bad  
.58 .42  
9%"]
    Node45 -- "no" --> Node48["Good  
.24 .76  
22%"]
    Node47 -- "volatile.acidity >= 0.35" --> Node49["Bad  
.58 .42  
9%"]
    Node47 -- "no" --> Node50["Good  
.24 .76  
22%"]
    Node49 -- "volatile.acidity >= 0.35" --> Node51["Bad  
.58 .42  
9%"]
    Node49 -- "no" --> Node52["Good  
.24 .76  
22%"]
    Node51 -- "volatile.acidity >= 0.35" --> Node53["Bad  
.58 .42  
9%"]
    Node51 -- "no" --> Node54["Good  
.24 .76  
22%"]
    Node53 -- "volatile.acidity >= 0.35" --> Node55["Bad  
.58 .42  
9%"]
    Node53 -- "no" --> Node56["Good  
.24 .76  
22%"]
    Node55 -- "volatile.acidity >= 0.35" --> Node57["Bad  
.58 .42  
9%"]
    Node55 -- "no" --> Node58["Good  
.24 .76  
22%"]
    Node57 -- "volatile.acidity >= 0.35" --> Node59["Bad  
.58 .42  
9%"]
    Node57 -- "no" --> Node60["Good  
.24 .76  
22%"]
    Node59 -- "volatile.acidity >= 0.35" --> Node61["Bad  
.58 .42  
9%"]
    Node59 -- "no" --> Node62["Good  
.24 .76  
22%"]
    Node61 -- "volatile.acidity >= 0.35" --> Node63["Bad  
.58 .42  
9%"]
    Node61 -- "no" --> Node64["Good  
.24 .76  
22%"]
    Node63 -- "volatile.acidity >= 0.35" --> Node65["Bad  
.58 .42  
9%"]
    Node63 -- "no" --> Node66["Good  
.24 .76  
22%"]
    Node65 -- "volatile.acidity >= 0.35" --> Node67["Bad  
.58 .42  
9%"]
    Node65 -- "no" --> Node68["Good  
.24 .76  
22%"]
    Node67 -- "volatile.acidity >= 0.35" --> Node69["Bad  
.58 .42  
9%"]
    Node67 -- "no" --> Node70["Good  
.24 .76  
22%"]
    Node69 -- "volatile.acidity >= 0.35" --> Node71["Bad  
.58 .42  
9%"]
    Node69 -- "no" --> Node72["Good  
.24 .76  
22%"]
    Node71 -- "volatile.acidity >= 0.35" --> Node73["Bad  
.58 .42  
9%"]
    Node71 -- "no" --> Node74["Good  
.24 .76  
22%"]
    Node73 -- "volatile.acidity >= 0.35" --> Node75["Bad  
.58 .42  
9%"]
    Node73 -- "no" --> Node76["Good  
.24 .76  
22%"]
    Node75 -- "volatile.acidity >= 0.35" --> Node77["Bad  
.58 .42  
9%"]
    Node75 -- "no" --> Node78["Good  
.24 .76  
22%"]
    Node77 -- "volatile.acidity >= 0.35" --> Node79["Bad  
.58 .42  
9%"]
    Node77 -- "no" --> Node80["Good  
.24 .76  
22%"]
    Node79 -- "volatile.acidity >= 0.35" --> Node81["Bad  
.58 .42  
9%"]
    Node79 -- "no" --> Node82["Good  
.24 .76  
22%"]
    Node81 -- "volatile.acidity >= 0.35" --> Node83["Bad  
.58 .42  
9%"]
    Node81 -- "no" --> Node84["Good  
.24 .76  
22%"]
    Node83 -- "volatile.acidity >= 0.35" --> Node85["Bad  
.58 .42  
9%"]
    Node83 -- "no" --> Node86["Good  
.24 .76  
22%"]
    Node85 -- "volatile.acidity >= 0.35" --> Node87["Bad  
.58 .42  
9%"]
    Node85 -- "no
```

Accuracy: 74.38% - Successfully classifies 3 out of 4 wines correctly

AUC: 0.7963 - Good discriminatory power between quality classes

Balanced Performance: Sensitivity (74.25%) and Specificity (74.51%) are nearly equal, indicating the model performs well on both Good and Bad wines

Accuracy: 74.38% - Successfully classifies 3 out of 4 wines correctly

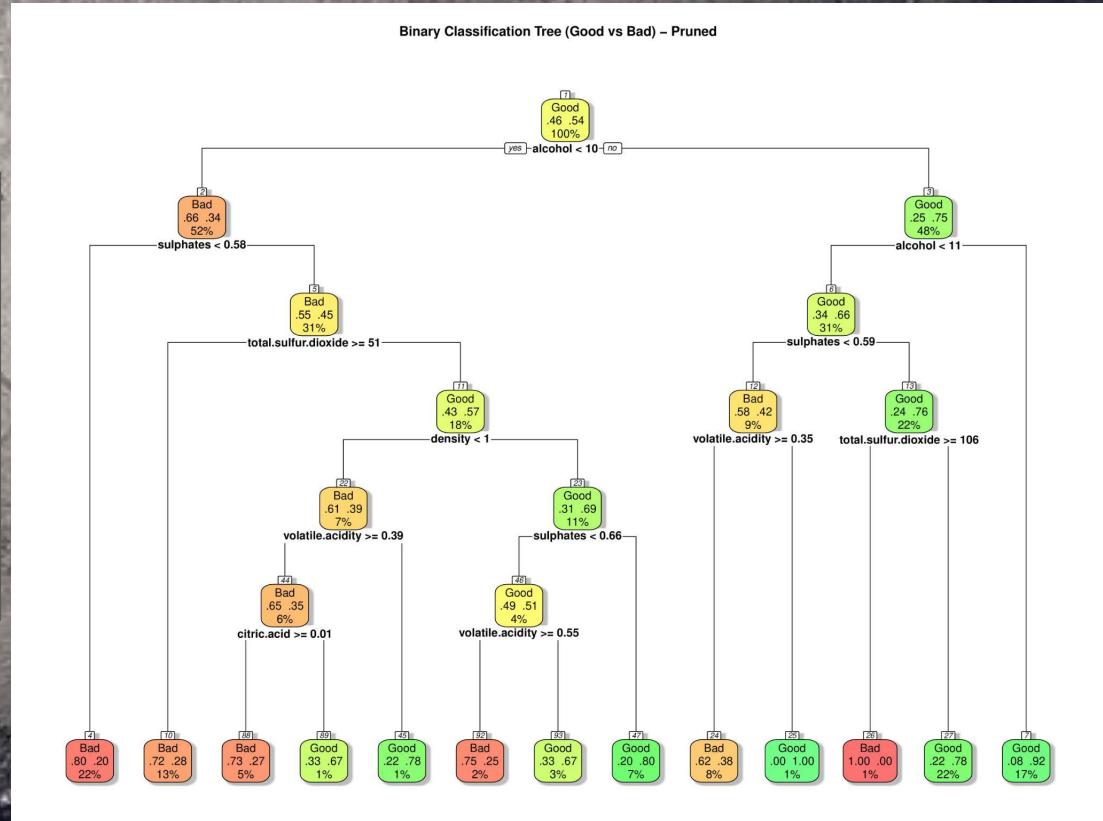
AUC: 0.7963 - Good discriminatory power between quality classes

Balanced Performance: Sensitivity (74.25%) and Specificity (74.51%) are nearly equal, indicating the model performs well on both Good and Bad wines

Accuracy: 74.38% - Successfully classifies 3 out of 4 wines correctly

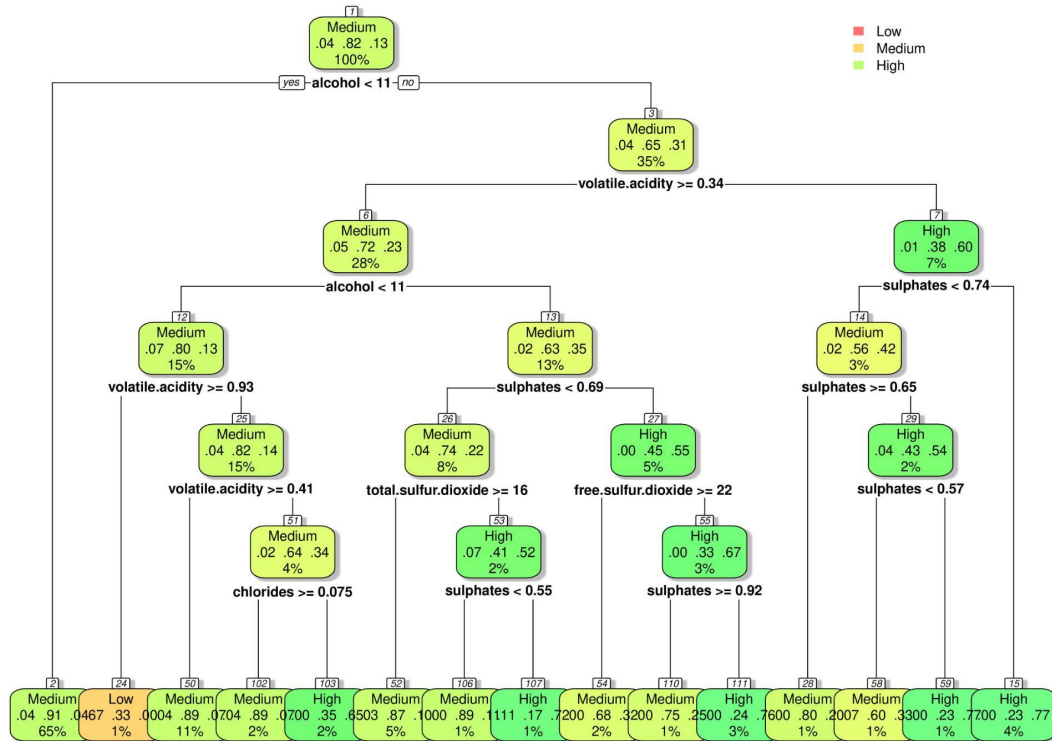
AUC: 0.7963 - Good discriminatory power between quality classes

Balanced Performance: Sensitivity (74.25%) and Specificity (74.51%) are nearly equal, indicating the model performs well on both Good and Bad wines



Multi-Class Classification Tree

Multi-class Classification Tree (Low/Medium/High) – Pruned



Overall Accuracy: 86.25% - Impressive performance

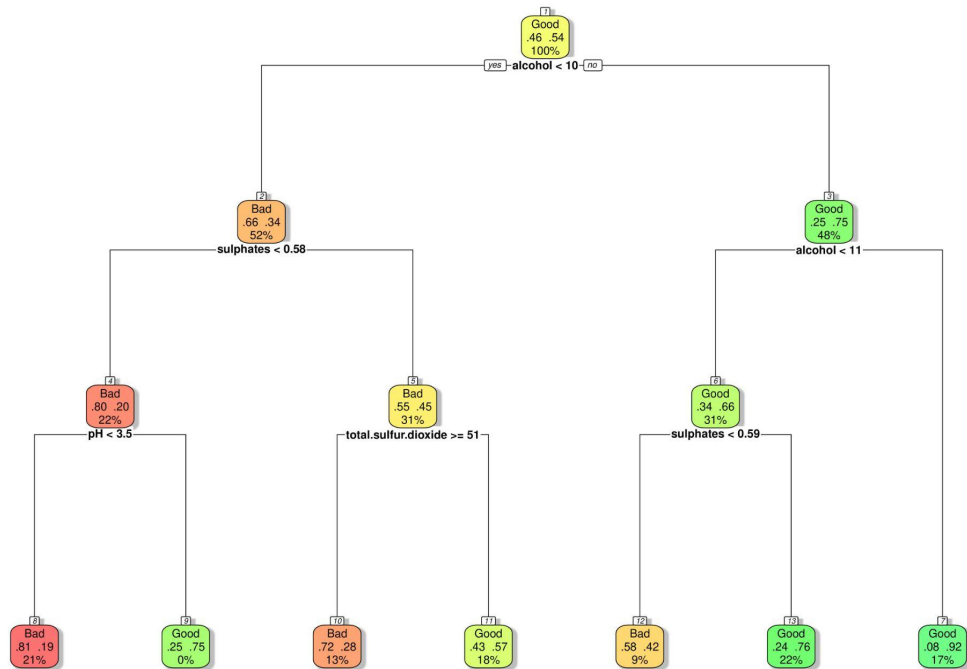
Strong at Medium Quality: 94.38% sensitivity for Medium wines (the majority class)

Higher-Alcohol Wines: alcohol content is the strongest overall predictor

Lower-alcohol wines: volatile acidity becomes the main driver

Hyperparameter-Tuned Binary Tree

Hyperparameter-Tuned Binary Tree
CP=0.001, minsplit=10, maxdepth=3



Tuned Model Performance: 72.5% accuracy

Slightly lower than pruned model (74.38%), but achieved better **Recall (79.04%)** at the expense of precision

Trade-off: The tuned model is more sensitive (catches more Good wines) but less precise (more false positives)

Gradient Boosting and Random Forest

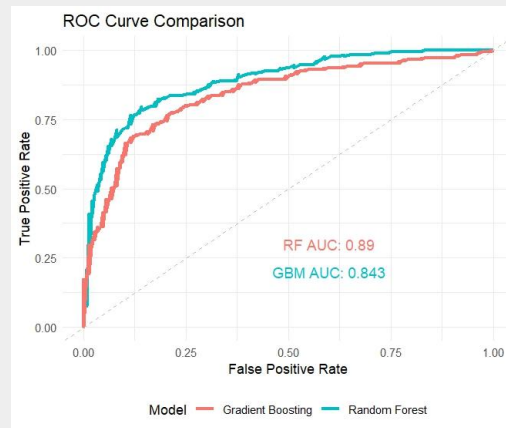
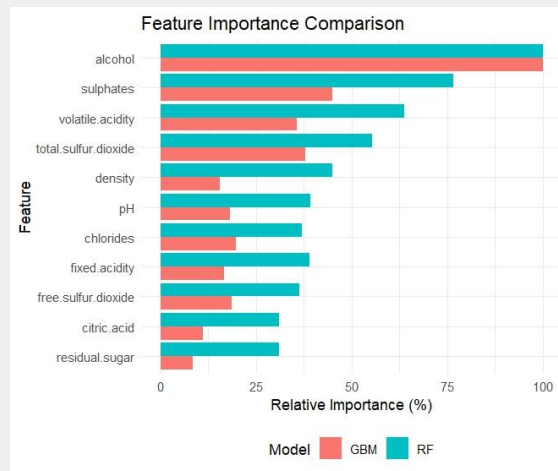
Consistent Top Predictors (Both Models Agree):

- **Alcohol** dominates in both models as the strongest predictor
- **Sulphates** ranks as the second most important feature
- **Volatile acidity** maintains strong importance

Random Forest emerges as the superior model with a test accuracy of **81.19%** compared to GBM's **77.74%**. The RF model demonstrates:

- Higher precision (**81.71% vs 79.41%**)
- **Better recall** (83.63% vs 78.95%)
- Superior F1-score (82.66% vs 79.18%)
- Notably stronger AUC (0.89 vs 0.843)

However, there's a concerning signal: RF shows **perfect training accuracy (100%)**, suggesting potential overfitting. This is why GBM's **training accuracy of 85.86%** indicates better generalization characteristics - it learns from errors of earlier trees, avoiding overfitting



Alcohol Threshold Analysis

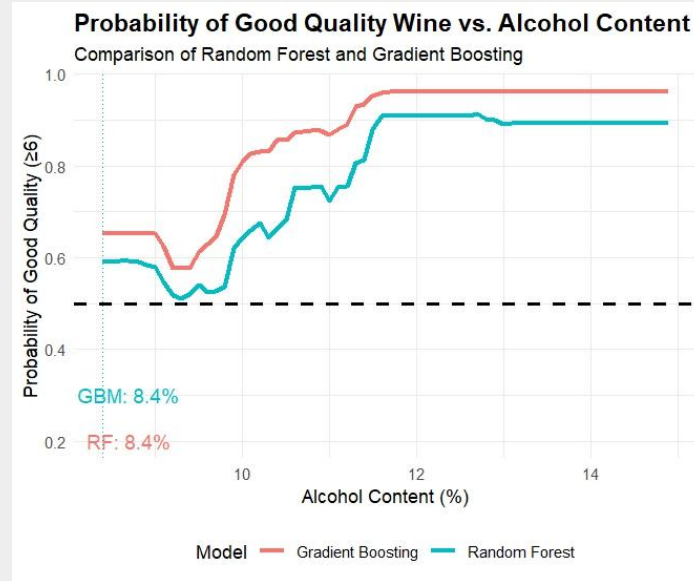
8.4% is the minimum alcohol content for both models

Gradient Boosting (Red Line):

- Starts at ~65% probability at 8.4%
- Rises smoothly to ~97% at 14%+
- Shows **continuous positive relationship**

Random Forest (Teal Line):

- Starts at ~60% probability at 8.4%
- More **step-like pattern** (characteristic of tree-based models)
- Plateaus at ~92% around 11-12%
- Shows some **non-monotonic behavior** (slight dips)



- **8.4% alcohol** = 60-65% good quality probability
- **10-11% alcohol** = 85-90% good quality probability (optimal target)
- **12%+ alcohol** = 90-95% good quality probability (marginal gains)

Main Findings

1. Can quality be predicted using only 3 chemical measurements?

✓ **3-feature model:** Yes—alcohol, volatile acidity, and sulphates together account for the model's predictive power. This represents a highly efficient feature set.

2. Is alcohol content alone sufficient, or are acidity measures necessary?

✗ **Alcohol alone is insufficient:** While alcohol is the dominant predictor, volatile acidity adds substantial predictive value. Using alcohol alone would sacrifice significant accuracy—the combined approach is necessary for optimal performance.

3. What minimum alcohol percentage gives >50% probability of good quality?

✓ **Threshold discovered:** "Threshold analysis reveals **8.4%** as the minimum alcohol level, but more importantly, it shows that quality depends on holistic chemical interactions. Winemakers should target **10-12%** alcohol while maintaining optimal chemical balance across all three key factors."

Conclusion

Key Takeaway: Wine quality prediction requires combination of alcohol content, volatile acidity, and sulphates - no single feature will suffice when predicting wine quality

Practical Value: Winemakers can target 10-12% alcohol threshold to improve quality ratings

Model Recommendation: Random Forest (**81.19% accuracy**) is optimal for quality control applications requiring maximum accuracy. Gradient Boosting (**77.74% accuracy**) is preferred for predictive applications for new, unseen wine batches

Impact: Provides actionable guidance for wine production optimization

Questions??



Why Use Hyper-Parameter Tuning?

Class imbalance: Dataset has more medium-quality wines than excellent/poor quality—tuning allows me to adjust the parameters like `min_samples_leaf` to prevent the model from creating tiny leaf nodes for rare classes, producing a balance of good and bad, not just medium quality

Reduces overfitting: Binary pruned tree didn't capture all the data; tuning the `CP`, `max_depth`, `minsplit`, allowed me to find optimal balance

Interpretability vs. Performance:
allows the tree to generalize better to new, unseen data rather than just memorizing the training set



Random Forest in Predicting Wine Quality

1. Bootstrap Aggregating (Bagging)

What Happens: Each of the **300 trees** in my Random Forest doesn't see all **1,280** training wines. Instead, each tree is trained on a random sample of **~1,280** wines drawn *with replacement* from the training set. This means:

- Some wines appear multiple times in a tree's training data
- Some wines don't appear at all (~37% are left out)
- Every tree sees a slightly different version of the wine quality story

2. Random Feature Selection (mtry)

What Happens: At each decision point in every tree, instead of considering all **11** chemical properties to make a split, the algorithm randomly selects only 4 features (**mtry=4**). So when a tree asks 'how should I divide these wines?', it might only consider:

- Alcohol, pH, density, chlorides (in one split)
- Volatile acidity, sulphates, citric acid, residual sugar (in another split)

This constraint forces trees to be different from each other

3. Majority Voting (Aggregation)

What Happens: "When predicting whether a new wine is 'Good' or 'Bad' quality:

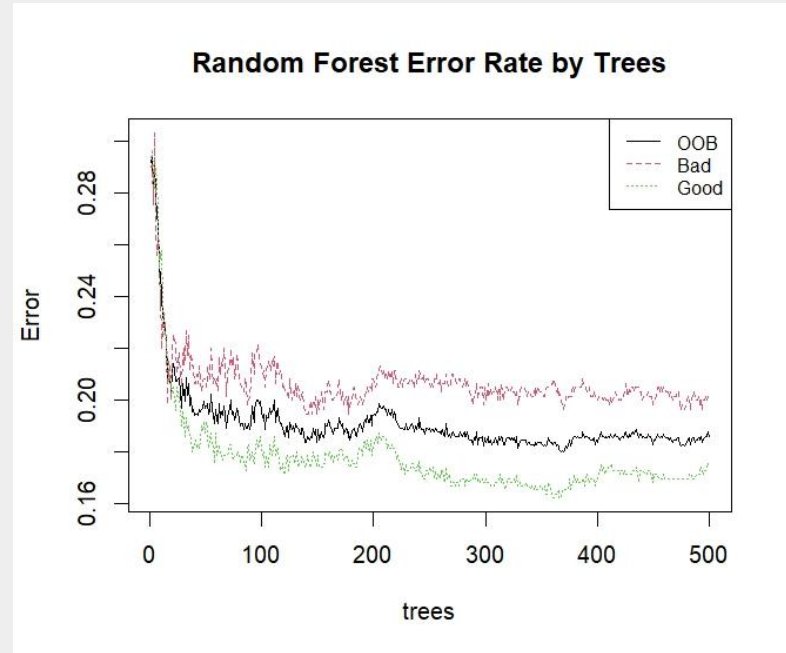
- Each of the 300 trees makes its own prediction independently
 - We count the votes: maybe **245** trees say 'Good', **55** say 'Bad'
 - The final prediction is 'Good' (majority wins)
 - The confidence/probability is **245/300 = 81.67%**

Random Forest Error Rate

- **mtry = 4** (number of variables randomly sampled at each split)
- **ntree = 300** (number of trees)
- **nodesize = 1** (minimum node size)
- **OOB Error = 17.81%**

The error rate plot shows:

- Initial high error (~**28%**) rapidly decreases in the first 50-100 trees
- Stabilization occurs around 150-200 trees
- Final OOB error stabilizes around **18%**
- Class-specific errors: Bad wines (~**20%**), Good wines (~**16-17%**)



Why pair Gradient Boosting with Random Forest?

Random Forest (Parallel Learning):

- Builds 300 independent trees simultaneously
- Each tree sees a random subset of features ($m_{try}=4$)
- Reduces variance through averaging
- **My finding:** Achieved **81.19%** test accuracy with strong feature separation

Gradient Boosting (Sequential Learning):

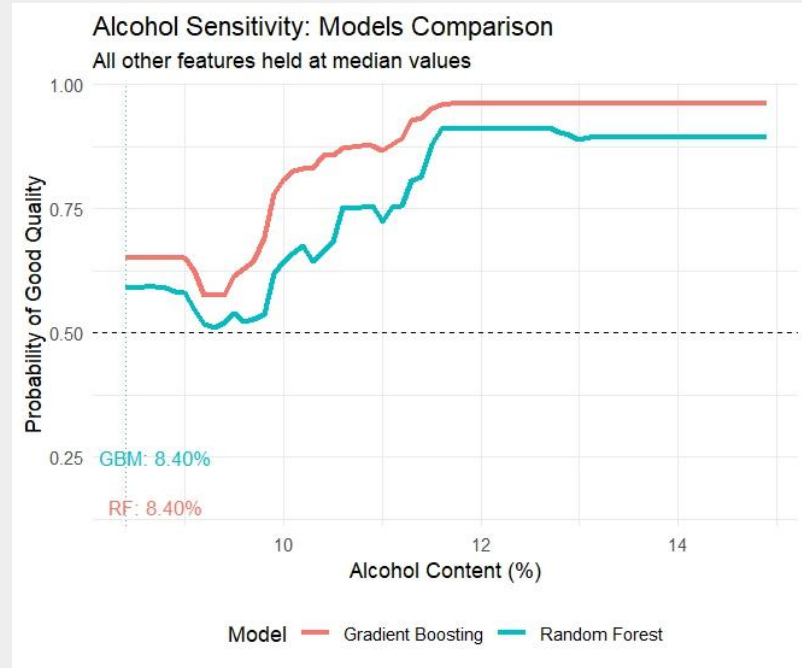
- Builds trees sequentially, each correcting previous errors
- Focuses on misclassified observations iteratively
- Reduces both bias AND variance
- Reduces overfitting

Alcohol Sensitivity Analysis

Median features baseline: 8.4% for both Random Forest and GBM

Wine makers should target **10-11%** alcohol as the "sweet spot" where probability of good quality increases most rapidly for typical wines

Alcohol content alone is NOT sufficient - acidity and other chemical measures are necessary



Alcohol Sensitivity By Wine Profile

- Wines with "good" chemical profiles (low volatile acidity, optimal sulphates) reach **50%** probability at just **8.4%** alcohol
- Wines with "bad" chemical profiles need nearly **10%** alcohol to reach the same probability
- Even at maximum alcohol (**15%**), bad-profile wines plateau around **72-75%** probability, while good-profile wines reach **95%+**

