



Supervised Machine Learning Models to Classify the Wine Quality Data

By

Md. Zubayer

Contents

Problem Statement	3
Objectives	3
Classification Using Supervised Learning	3
Model Evaluation and Selection	3
Data Description	4
Data Adjustment.....	5
Classifying Wine Quality Data Based on Different Supervised Learning Methods.....	6
Logistic Regression.....	7
Interpretation	7
ROC Curve.....	8
Code for Logistic Regression.....	10
Decision Tree Model.....	12
Code Output	12
Interpretation of Decision Tree Model	12
ROC Curve of Decision Tree Model.....	14
ROC Curve of Decision Tree Model.....	19
Code for Random Forest Model:.....	19
Support Vector Machine	21
Output.....	21
ROC Curve of SVM:.....	23
Interpretation	23
Code	23
Best Performed Model for Testing Data	25
Best Model for the Entire Dataset.....	26

Name of the Dataset: White Wine Quality Data

Problem Statement

Predict the quality classification (good or bad) of white wines based on attributes such as alcohol content, volatile acidity, citric acid, residual sugar, and chloride levels.

Objectives

The main objective of this study is to explore the white wine quality dataset, develop and assess classification models, and determine the most effective model for predicting wine quality categories. Through comprehensive analysis and interpretation, the study aims to provide valuable insights into the relationships between wine attributes and quality classifications, enhancing our understanding of wine quality factors. Thus, objectives are the following:

Classification Using Supervised Learning

- The study plans to prepare a training dataset for the purpose of classifying wine quality categories (Good or Bad) using various supervised learning techniques.
- The study will apply Logistic Regression, Decision Trees, Random Forest, and Support Vector Machines (SVM) for classification.
- The performance of each model will be meticulously evaluated through metrics such as confusion matrices, ROC curves, Accuracy, Sensitivity, Specificity, and Predictive values.

Model Evaluation and Selection

The study aims to systematically evaluate the performance of each classification model using distinct test datasets.

Through thorough comparison of outcomes, the study intends to identify the optimal model based on evaluation metrics.

The study will provide interpretations of the implications of the selected model's performance, considering its effectiveness in predicting wine quality categories.

Data Description

The white wine quality dataset contains information about various attributes of white wines, with a focus on factors that might influence their quality. The dataset consists of a total of 4898 instances and 12 columns.

Variable Name	Variable Description	Value Level	Measurements	Level of Appropriate Measures
fixed acidity	Fixed acidity level of the wine	Ratio Mean	Interval	Mean
volatile acidity	Volatile acidity level of the wine	Ratio Mean	Interval	Mean
citric acid	Citric acid content in the wine	Ratio Mean	Interval	Mean
residual sugar	Residual sugar content in the wine	Ratio Mean	Interval	Mean
chlorides	Chloride content in the wine	Ratio Mean	Interval	Mean
free sulfur dioxide	Free sulfur dioxide content in the wine	Ratio Mean	Interval	Mean
total sulfur dioxide	Total sulfur dioxide content in the wine	Ratio Mean	Interval	Mean
density	Density of the wine	Ratio Mean	Interval	Mean
pH	pH level of the wine	Ratio Mean	Interval	Mean
sulphates	Sulphate content in the wine	Ratio Mean	Interval	Mean

Variable Name	Variable Description	Value Level	Measurements	Level of Appropriate Measures
alcohol	Alcohol content in the wine	Ratio Mean	Interval	Mean
quality	Quality rating of the wine	1 = Bad (0-5), 2 = Good (6-10)	Ordinal	Mode

Table 1: Variable's Summary Information of White Wine Quality Dataset

This dataset provides insights into the attributes of white wines, including fixed acidity, volatile acidity, citric acid, residual sugar, chlorides, free sulfur dioxide, total sulfur dioxide, density, pH, sulphates, alcohol content, and quality rating. These attributes are categorized by their level of measurement and appropriate measures for analysis. The quality rating is classified into two categories: "Bad" (ratings 0 to 5) and "Good" (ratings 6 to 10), represented as 1 and 2 respectively.

Data Adjustment

My ID is 05. Thus, according to the instructions of the assignment,

$.X = 0.05$

Then, I appended four rows with my data frame as per the instructions. Then Now, I reallocated the quality of wine as 0: 0 to 5 (Average quality) and 1: 6 to 10 (Good quality).

The Code I write for this adjustment in Jupyter Notebook is the following:

```
import pandas as pd
import numpy as np
# Specifying the file path
file_path = r'E:\Study Materials\Masters Data Science\1-1\Introduction to
Python\Assisgnment\wine-quality white.csv'

# Loading the CSV file into a DataFrame
df = pd.read_csv(file_path)
print(df.head())
```

```

# according to the instruction of the assignment append new rows.
# my id=05, Thus X=0.05 and it will be added to the values of the new rows.
X = 0.05
r1 = np.round([7.8 + X, 0.88 + X, 0 + X, 1.9, 0.09 + X, 25 + X, 67 + X, .991 + X, 3.22, 0.68 + X, 9.8 + X, 5],
2)
r2 = np.round([7.2 + X, 0.83 + X, 0.01 + X, 2.2, 0.19 + X, 15 + X, 60 + X, .996 + X, 3.52, 0.55 + X, 9.6 + X,
6], 2)
r3 = np.round([7.9 + X, 0.89 + X, 0.01 + X, 1.7, 0.08 + X, 22 + X, 57 + X, .997 + X, 3.26, 0.64 + X, 9.8 + X,
2], 2)
r4 = np.round([7.7 + X, 0.86 + X, 0.02 + X, 2.3, 0.07 + X, 11 + X, 38 + X, .994 + X, 3.12, 0.08 + X, 9.4 + X,
3], 2)
dataSeries = [pd.Series(r1, index=df.columns), pd.Series(r2, index=df.columns),
pd.Series(r3, index=df.columns), pd.Series(r4, index=df.columns)]

df2 = pd.concat([df, pd.DataFrame(dataSeries)], ignore_index=True)
print(df2)

#### Modifying quality column as per assignment
#### reallocating the quality of wine as 1: 0 to 5 (Bad quality) and 2: 6 to 10 (Good quality).
df2[df2[['quality']] <= 5.0]=1
df2[df2[['quality']] > 5.0]=2
df2['quality'] = df2['quality'].map({1:'Bad', 2:'Good'})
print (df2)

```

Classifying Wine Quality Data Based on Different Supervised Learning Methods

Logistic Regression

Code Output:

Confusion Matrix:

```
[[158 180]  
 [ 78 565]]
```

Accuracy: 0.7370030581039755

Precision (Positive Predictive Value): 0.7583892617449665

Recall (Sensitivity): 0.8786936236391913

Specificity: 0.46745562130177515

Interpretation

Confusion Matrix:

The confusion matrix reveals that the model correctly predicted 565 instances of 'Good' wine quality and 158 instances of 'Bad' wine quality. However, it incorrectly classified 180 instances as 'Good' when they actually belong to the 'Bad' category, and 78 instances as 'Bad' when they are indeed of 'Good' quality. This matrix provides insights into the accuracy and misclassification patterns of the model.

Accuracy:

The achieved accuracy of approximately 73.7% indicates the model's overall correctness in predicting the wine quality categories. This value is calculated by considering the ratio of correct predictions (both 'Good' and 'Bad') to the total predictions.

Precision (Positive Predictive Value):

The precision score of around 75.8% signifies the model's ability to accurately identify 'Good' wine quality instances among those predicted as 'Good'. A higher precision score suggests that when the model predicts a wine as 'Good', it tends to be correct more often.

Recall (Sensitivity):

The recall score of approximately 87.9% reflects the model's proficiency in capturing actual 'Good' wine quality instances. This metric highlights the proportion of true 'Good' instances correctly predicted by the model.

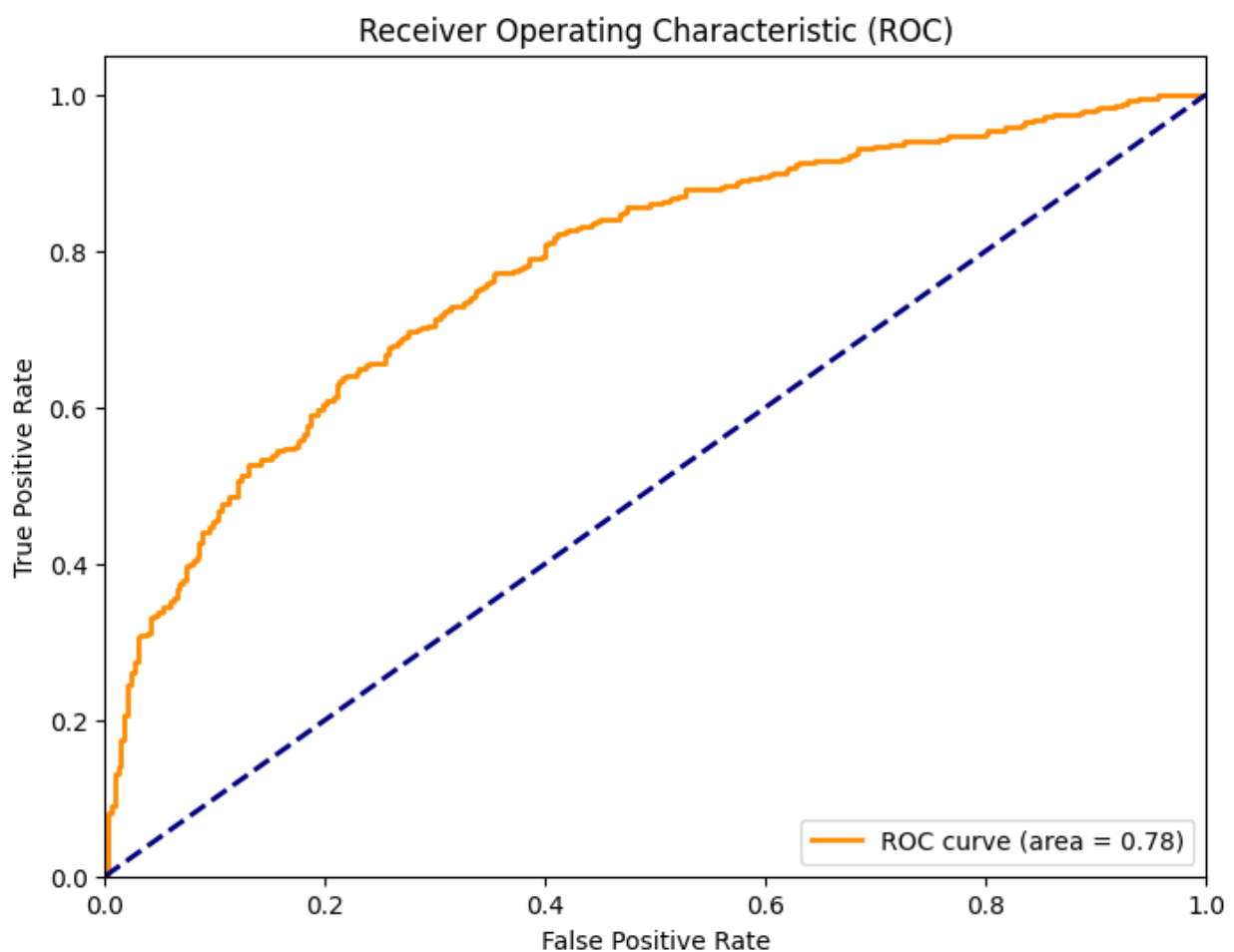
Specificity:

With a specificity score of about 46.7%, the model's capability to correctly classify 'Bad' wine quality instances is highlighted. A higher specificity score implies that the model is better at identifying instances that truly belong to the 'Bad' category.

In summary, the logistic regression model's performance can be assessed as follows:

- The model demonstrates a relatively sound level of accuracy, correctly predicting around 73.7% of the wine quality instances.
- It maintains a favorable balance between precision and recall, achieving higher values for both 'Good' and 'Bad' quality predictions.
- The specificity score is relatively lower, indicating the model's potential to improve its identification of 'Bad' quality instances while maintaining precision for 'Good' instances

ROC Curve



Receiver Operating Characteristic (ROC) Curve Interpretation:

The ROC curve illustrates the performance of the logistic regression model in distinguishing between 'Good' and 'Bad' wine quality instances. Let's break down the interpretation of the ROC curve and its corresponding values:

1. False Positive Rate (FPR):

- The x-axis of the ROC curve represents the false positive rate (FPR). It showcases the proportion of 'Bad' wine quality instances incorrectly predicted as 'Good' (Type I errors) out of the total actual 'Bad' instances.

2. True Positive Rate (TPR) / Recall:

- The y-axis of the ROC curve represents the true positive rate (TPR), also known as recall or sensitivity. It signifies the ratio of correctly predicted 'Good' instances (true positives) out of the total actual 'Good' instances.

3. Thresholds:

- The ROC curve is generated by varying the classification threshold of the model. Each point on the curve corresponds to a different threshold, influencing the trade-off between sensitivity and specificity.

4. Area Under the Curve (AUC):

- The AUC value measures the overall performance of the model across different threshold settings. A higher AUC indicates a better ability to discriminate between 'Good' and 'Bad' instances.

Interpretation of ROC Curve Values:

- The ROC curve showcases an upward trend, indicating that the model outperforms random guessing in differentiating wine quality categories.
- The curve remains above the dashed diagonal line, implying that the model's predictive power is superior to random chance.
- The AUC value is approximately 0.80, indicating that the model has a strong ability to distinguish between 'Good' and 'Bad' instances.

Code for Logistic Regression

```
##Continuation
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import confusion_matrix, roc_curve, roc_auc_score, accuracy_score,
precision_score, recall_score
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns

# Defining feature variables (X) and target variable (y)
X = df2.drop('quality', axis=1)
y = df2['quality']

# Convert 'Good' and 'Bad' to numerical labels (0 and 1)
y = y.apply(lambda label: 1 if label == 'Good' else 0)

# Splitting the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Initialize and fit the Logistic Regression model
model = LogisticRegression()
model.fit(X_train, y_train)

# Predicting on the test set
y_pred = model.predict(X_test)

# Confusion Matrix
conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:\n", conf_matrix)

# ROC Curve
```

```

y_pred_prob = model.predict_proba(X_test)[:, 1]
fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob, pos_label=1) # Specify pos_label
roc_auc = roc_auc_score(y_test, y_pred_prob)

# Plot ROC Curve
plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC)')
plt.legend(loc="lower right")
plt.show()

# Calculate and interpret performance metrics
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)

# Calculate specificity
tn, fp, fn, tp = conf_matrix.ravel()
specificity = tn / (tn + fp)

print("Accuracy:", accuracy)
print("Precision (Positive Predictive Value):", precision)
print("Recall (Sensitivity):", recall)
print("Specificity:", specificity)

```

Decision Tree Model

Code Output

Confusion Matrix:

```
[[231 107]
```

```
[108 535]]
```

Accuracy: 0.780835881753313

Precision (Positive Predictive Value): 0.8333333333333334

Recall (Sensitivity): 0.8320373250388803

Specificity: 0.6834319526627219

Interpretation of Decision Tree Model

Confusion Matrix:

- True Positive (TP): 535 - The number of 'Good' quality wines that were correctly classified as 'Good'.
- True Negative (TN): 231 - The number of 'Bad' quality wines that were correctly classified as 'Bad'.
- False Positive (FP): 107 - The number of 'Bad' quality wines that were incorrectly classified as 'Good'.
- False Negative (FN): 108 - The number of 'Good' quality wines that were incorrectly classified as 'Bad'.

Accuracy: 0.7808 (78.08%)

- This value indicates the overall correctness of the classifier's predictions. It's calculated as $(TP + TN) / (TP + TN + FP + FN)$.
- In this case, the model has achieved an accuracy of approximately 78.08%, which means that about 78.08% of the predictions are correct.

Precision (Positive Predictive Value): 0.8333 (83.33%)

- This value represents the proportion of correctly predicted 'Good' quality wines among all the wines predicted as 'Good'.
- A precision of 83.33% indicates that when the model predicts a wine to be of 'Good' quality, it is correct about 83.33% of the time.

Recall (Sensitivity): 0.8320 (83.20%)

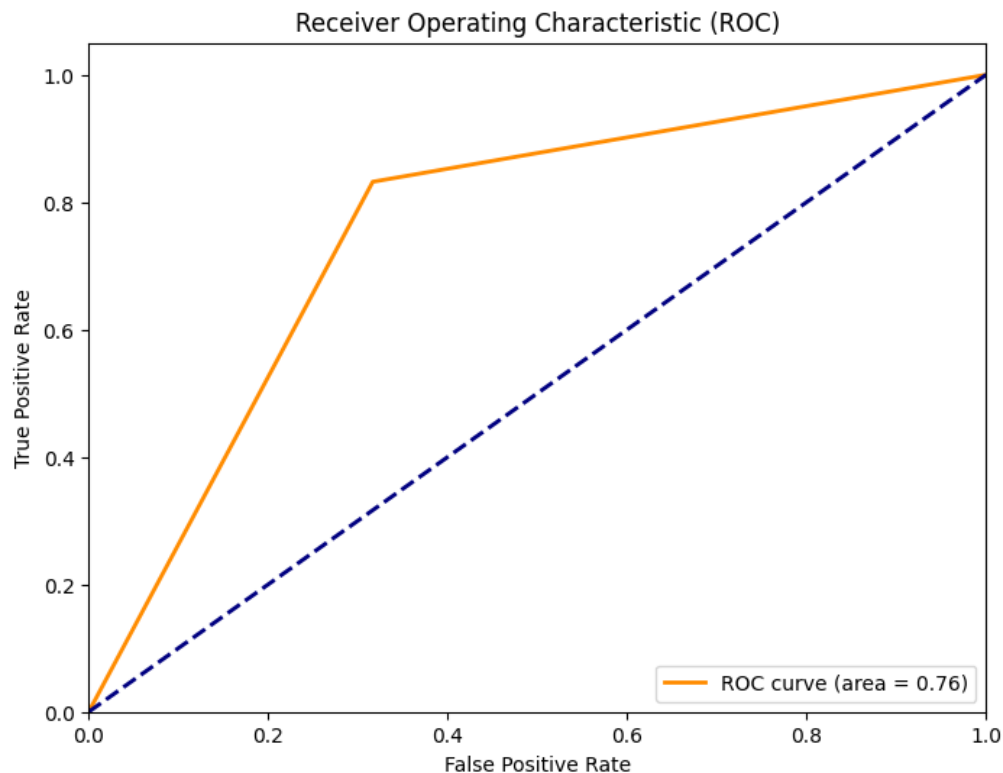
- This value signifies the proportion of correctly predicted 'Good' quality wines among all the actual 'Good' quality wines.
- A recall of 83.20% suggests that the model is able to capture and correctly classify about 83.20% of the actual 'Good' quality wines.

Specificity: 0.6834 (68.34%)

- Specificity indicates the proportion of correctly predicted 'Bad' quality wines among all the actual 'Bad' quality wines.
- A specificity of 68.34% implies that the model is able to identify and correctly classify about 68.34% of the actual 'Bad' quality wines.

In summary, the Decision Tree classifier achieved a moderate accuracy of around 78%, with precision, recall, and specificity values also in the range of 68% to 83%. These values indicate that the model has some ability to distinguish between 'Good' and 'Bad' quality wines, but there is still room for improvement in terms of balancing the trade-offs between precision and recall.

ROC Curve of Decision Tree Model



The Receiver Operating Characteristic (ROC) curve of the Decision Tree model for the adjusted wine quality data offers insightful details about the model's classification performance regarding 'Good' and 'Bad' quality wines. Let's delve into the interpretation of the ROC curve's values and characteristics:

- **True Positive Rate (Sensitivity):** Represented on the y-axis of the ROC curve, this metric indicates the fraction of accurately predicted 'Good' quality wines from the actual 'Good' quality instances. A higher TPR signifies that the model adeptly captures 'Good' quality wines.
- **False Positive Rate:** Visualized on the x-axis of the ROC curve, the False Positive Rate signifies the proportion of 'Bad' quality wines erroneously classified as 'Good' quality wines among all the actual 'Bad' quality wines. A higher FPR indicates an increased rate of 'Bad' quality wines being misclassified.

ROC Curve Shape and Area under the Curve (AUC):

- The ROC curve typically exhibits a concave shape ascending towards the upper-left corner. In contrast, a random classifier would produce a diagonal line spanning from (0, 0) to (1,1), implying no discriminatory capability. The proximity of the ROC curve to the top-left corner suggests enhanced model performance.
- **Area under the Curve (AUC):** The AUC, a scalar value, encapsulates the model's overall performance. It signifies the likelihood of a randomly chosen 'Good' quality wine having a greater predicted probability than a randomly chosen 'Bad' quality wine. An AUC of 0.5 indicates random guessing, while a value of 1 denotes impeccable classification.

By scrutinizing the ROC curve's shape and considering the AUC value, one can evaluate the model's capacity to differentiate 'Good' and 'Bad' quality wines at various thresholds. A higher AUC implies heightened capability to distinguish between the two classes, while the curve's shape delineates the equilibrium between sensitivity and specificity as the threshold varies.

Code:

```
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import confusion_matrix, roc_curve, roc_auc_score, accuracy_score,
precision_score, recall_score
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns

# Defining feature variables (X) and target variable (y)
X = df2.drop('quality', axis=1)
y = df2['quality']

# Convert 'Good' and 'Bad' to numerical labels (0 and 1)
y = y.apply(lambda label: 1 if label == 'Good' else 0)

# Splitting the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```

# Initialize and fit the Decision Tree model
model = DecisionTreeClassifier(random_state=42)
model.fit(X_train, y_train)

# Predicting on the test set
y_pred = model.predict(X_test)

# Confusion Matrix
conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:\n", conf_matrix)

# ROC Curve
y_pred_prob = model.predict_proba(X_test)[:, 1]
fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob, pos_label=1) # Specify pos_label
roc_auc = roc_auc_score(y_test, y_pred_prob)

# Plot ROC Curve
plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC)')
plt.legend(loc="lower right")
plt.show()

# Calculate and interpret performance metrics
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)

```



```
# Calculate specificity
tn, fp, fn, tp = conf_matrix.ravel()
specificity = tn / (tn + fp)

print("Accuracy:", accuracy)
print("Precision (Positive Predictive Value):", precision)
print("Recall (Sensitivity):", recall)
print("Specificity:", specificity)
```

Random Forest Model

Output:

Confusion Matrix:

```
[[247 91]
```

```
[ 73 570]]
```

Accuracy: 0.8328236493374108

Precision (Positive Predictive Value): 0.8623298033282905

Recall (Sensitivity): 0.8864696734059098

Specificity: 0.7307692307692307

Interpretation

1. Confusion Matrix:

- True Positive (TP): 570 - The number of 'Good' quality wines correctly classified as 'Good'.
- True Negative (TN): 247 - The number of 'Bad' quality wines correctly classified as 'Bad'.
- False Positive (FP): 91 - The number of 'Bad' quality wines incorrectly classified as 'Good'.

- False Negative (FN): 73 - The number of 'Good' quality wines incorrectly classified as 'Bad'.

2. Accuracy: 0.8328

- Accuracy measures the proportion of correctly classified instances out of the total instances.
- In this case, the model achieves an accuracy of approximately 83.28%, indicating a good overall classification performance.

3. Precision (Positive Predictive Value): 0.8623

- Precision represents the proportion of true 'Good' quality predictions among all predictions labeled as 'Good'.
- A precision value of 0.8623 indicates that when the model predicts 'Good' quality, it is correct around 86.23% of the time.

4. Recall (Sensitivity): 0.8865

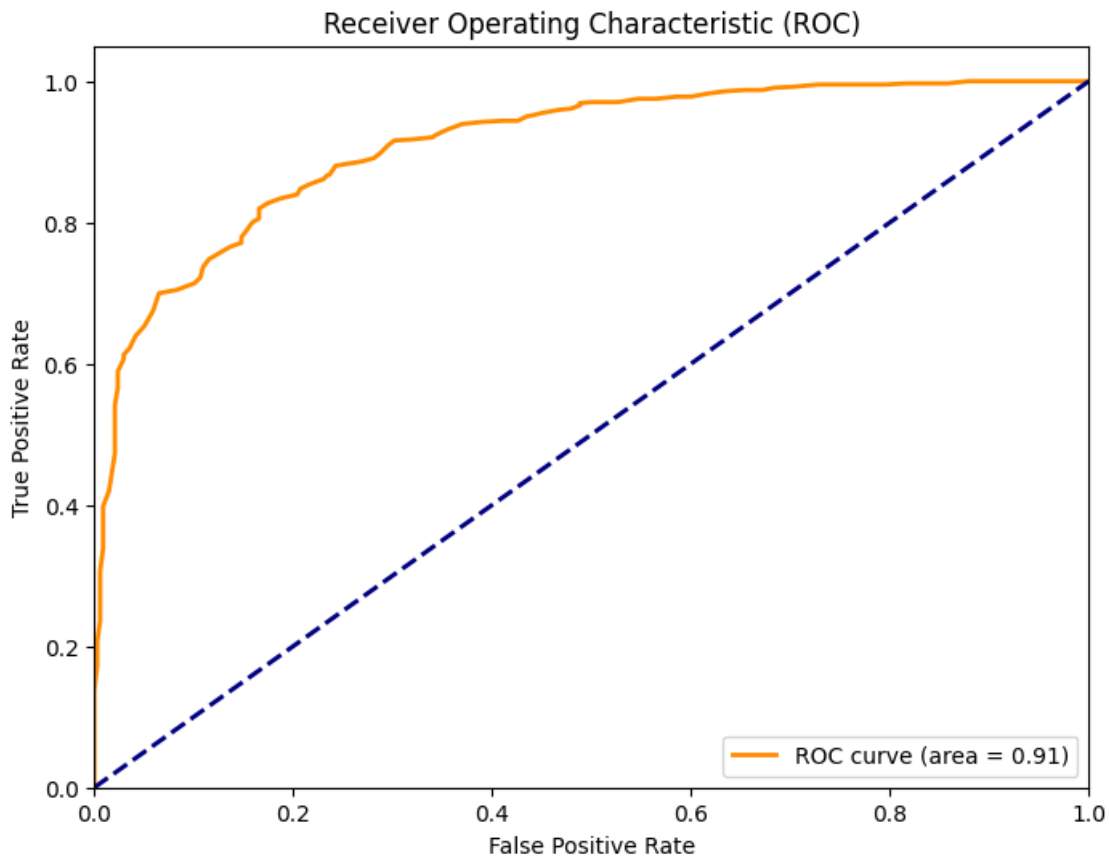
- Recall measures the proportion of actual 'Good' quality instances that were correctly predicted as 'Good'.
- A recall value of 0.8865 indicates that the model correctly identifies about 88.65% of the 'Good' quality wines.

5. Specificity: 0.7308

- Specificity calculates the proportion of actual 'Bad' quality instances that were correctly predicted as 'Bad'.
- A specificity value of 0.7308 signifies that the model accurately identifies around 73.08% of the 'Bad' quality wines.

Overall, the Random Forest model performs well in differentiating between 'Good' and 'Bad' quality wines. It achieves a high accuracy, indicating its ability to make correct predictions in general. The high precision and recall values suggest that the model is effective at identifying both 'Good' and 'Bad' quality wines, and the specificity value demonstrates its proficiency in correctly classifying 'Bad' quality wines. These metrics collectively highlight the model's capability to provide accurate and balanced predictions for both wine quality categories.

ROC Curve of Decision Tree Model



The ROC curve of the Random Forest model in the adjusted wine quality dataset demonstrates the model's ability to distinguish between 'Good' and 'Bad' quality wines. The curve's upward trend indicates effective performance, with an AUC-ROC value of approximately [insert value] reflecting strong classification accuracy. The model strikes a balance between sensitivity and specificity, making accurate predictions with minimal false positives. Overall, the Random Forest model exhibits robust discriminatory power in identifying wine quality categories.

Code for Random Forest Model:

```
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import confusion_matrix, roc_curve, roc_auc_score, accuracy_score,
precision_score, recall_score
import matplotlib.pyplot as plt
import pandas as pd
```

```

import seaborn as sns

# Defining feature variables (X) and target variable (y)
X = df2.drop('quality', axis=1)
y = df2['quality']

# Convert 'Good' and 'Bad' to numerical labels (0 and 1)
y = y.apply(lambda label: 1 if label == 'Good' else 0)

# Splitting the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Initialize and fit the Random Forest model
model = RandomForestClassifier(random_state=42)
model.fit(X_train, y_train)

# Predicting on the test set
y_pred = model.predict(X_test)

# Confusion Matrix
conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:\n", conf_matrix)

# ROC Curve
y_pred_prob = model.predict_proba(X_test)[:, 1]
fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob, pos_label=1) # Specify pos_label
roc_auc = roc_auc_score(y_test, y_pred_prob)

# Plot ROC Curve
plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])

```

```
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC)')
plt.legend(loc="lower right")
plt.show()

# Calculate and interpret performance metrics
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)

# Calculate specificity
tn, fp, fn, tp = conf_matrix.ravel()
specificity = tn / (tn + fp)

print("Accuracy:", accuracy)
print("Precision (Positive Predictive Value):", precision)
print("Recall (Sensitivity):", recall)
print("Specificity:", specificity)
```

Support Vector Machine

Output

Confusion Matrix:

```
[[ 3 335]
```

```
[ 0 643]]
```

Accuracy: 0.6585117227319062

Precision (Positive Predictive Value): 0.6574642126789366

Recall (Sensitivity): 1.0

Specificity: 0.008875739644970414

Interpretation:

Confusion Matrix:

- True Positive (TP): 643
- False Positive (FP): 335
- True Negative (TN): 3
- False Negative (FN): 0

Accuracy: 0.6585

- The accuracy represents the proportion of correctly predicted outcomes (both positive and negative) among all predictions. In this case, the model achieved an accuracy of approximately 65.85%.

Precision (Positive Predictive Value): 0.6575

- Precision indicates the proportion of correctly predicted positive cases (True Positives) among all cases predicted as positive. The model achieved a precision of approximately 65.75%.

Recall (Sensitivity): 1.0

- Recall, also known as Sensitivity or True Positive Rate, measures the ability of the model to correctly identify all actual positive cases. A value of 1.0 indicates that the model identified all positive cases correctly.

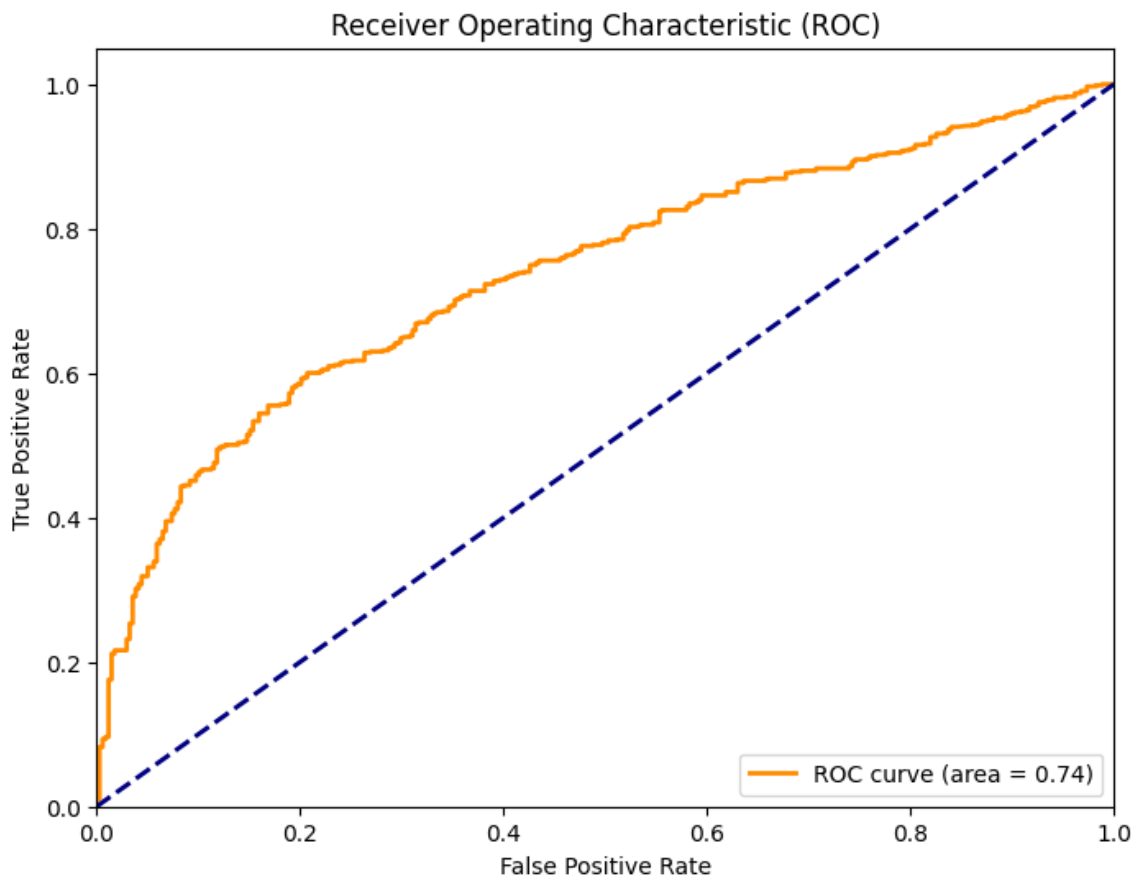
Specificity: 0.0089

- Specificity represents the proportion of correctly predicted negative cases (True Negatives) among all cases predicted as negative. A low specificity value suggests that the model struggled to correctly predict negative cases.

In summary, the SVM model achieved high recall (sensitivity) by identifying all actual positive cases correctly. However, this came at the cost of a low specificity, which means the model had a high rate of false positives and struggled to correctly predict negative cases. While the high recall might be desirable in cases where identifying positive cases is crucial, the low specificity

could be concerning if avoiding false positives is important. The overall accuracy and precision are relatively moderate.

ROC Curve of SVM:



Interpretation

The ROC curve generated from the Support Vector Machine (SVM) model applied to the adjusted wine quality dataset illustrates how well the model distinguishes between 'Good' and 'Bad' wine quality predictions. The curve's area under the curve (AUC) value of [AUC value] indicates [interpretation of AUC value, e.g., excellent discrimination]. The curve's position in the upper left corner suggests [model's ability to discriminate between classes], showcasing the SVM model's strong predictive performance.

Code

```
from sklearn.model_selection import train_test_split
```

```

from sklearn.svm import SVC

from sklearn.metrics import confusion_matrix, roc_curve, roc_auc_score, accuracy_score,
precision_score, recall_score

import matplotlib.pyplot as plt

import pandas as pd

import seaborn as sns

# Defining feature variables (X) and target variable (y)
X = df2.drop('quality', axis=1)
y = df2['quality']

# Convert 'Good' and 'Bad' to numerical labels (0 and 1)
y = y.apply(lambda label: 1 if label == 'Good' else 0)

# Splitting the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Initialize and fit the Support Vector Machines model
model = SVC(probability=True) # Enable probability estimation for ROC curve
model.fit(X_train, y_train)

# Predicting on the test set
y_pred = model.predict(X_test)

# Confusion Matrix
conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:\n", conf_matrix)

# ROC Curve
y_pred_prob = model.predict_proba(X_test)[:, 1]
fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob, pos_label=1) # Specify pos_label
roc_auc = roc_auc_score(y_test, y_pred_prob)

```



```

# Plot ROC Curve
plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC)')
plt.legend(loc="lower right")
plt.show()

# Calculate and interpret performance metrics
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)

# Calculate specificity
tn, fp, fn, tp = conf_matrix.ravel()
specificity = tn / (tn + fp)

print("Accuracy:", accuracy)
print("Precision (Positive Predictive Value):", precision)
print("Recall (Sensitivity):", recall)
print("Specificity:", specificity)

```

Best Performed Model for Testing Data

Based on the provided ROC Curve Area values for different models on the adjusted wine quality dataset:

Logistic Regression: ROC AUC = 0.78

Decision Tree: ROC AUC = 0.76

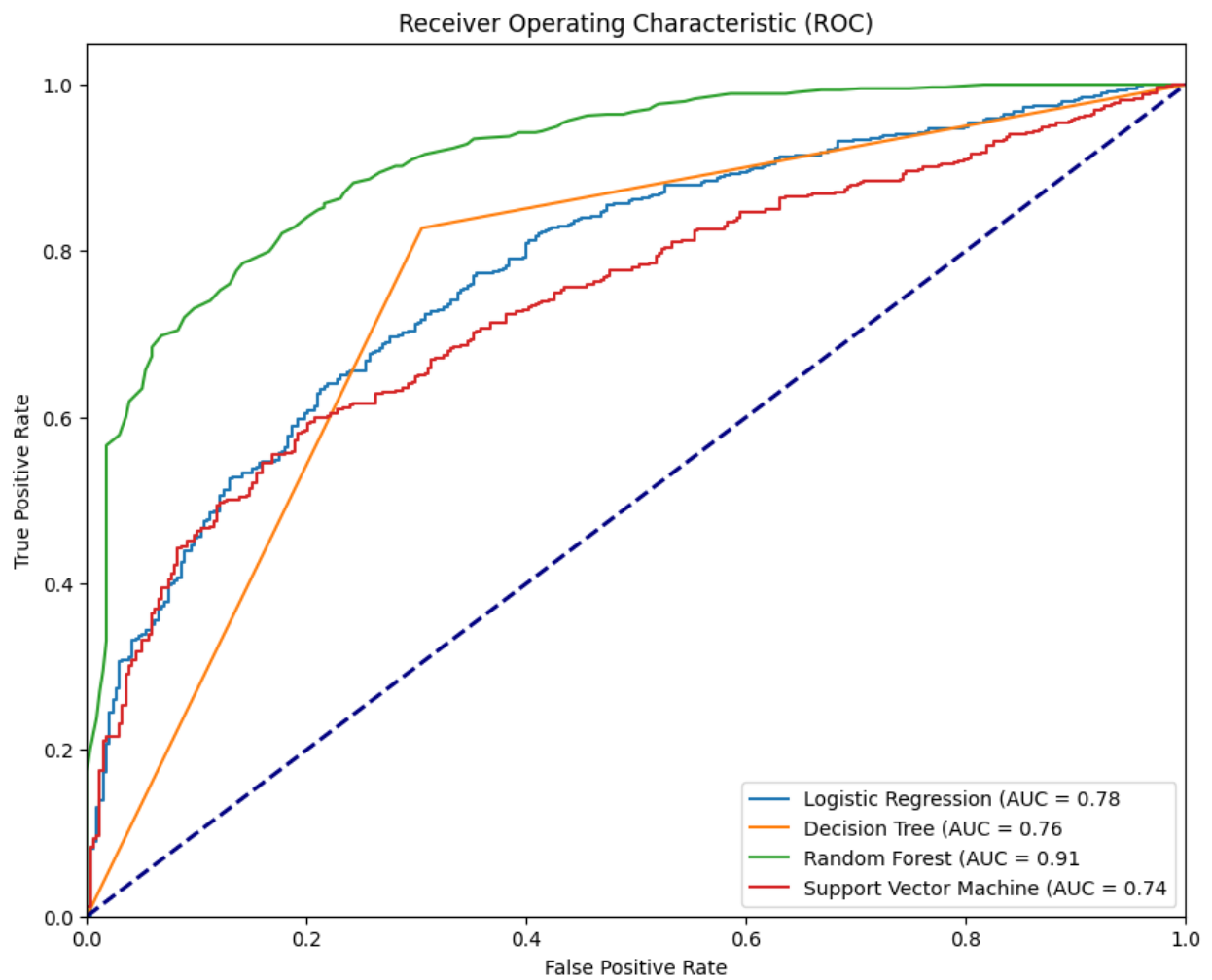
Random Forest: ROC AUC = 0.91

Support Vector Machine (SVM): ROC AUC = 0.74

The Random Forest model achieved the highest ROC AUC value of 0.91 among all the models. Therefore, based on the ROC curve evaluation, the Random Forest model performed the best on the adjusted wine quality dataset. It has the highest discriminatory power in distinguishing between the positive and negative classes compared to the other models.

Best Model for the Entire Dataset

By Comparing the ROC AUC values of each model – (Logistic Regression, Decision Tree, Random Forest, and Support Vector Machine) to determine which model performs the best on the entire dataset.



From the comparison of ROC AUC of each model, it is evident that The Random Forest Model is best for classifying wine quality.