

TASK1: Project Description

The dataset is related to red and white variants of the Portuguese "Vinho Verde" wine. Due to privacy and logistic issues, only physicochemical (inputs) and sensory (the output) variables are available (e.g. there is no data about grape types, wine brand, wine selling price, etc.).

This dataset can be viewed as classification task. The classes are ordered and not balanced (e.g. there are many more normal wines than excellent or poor ones). Also, we are not sure if all input variables are relevant. So it could be interesting to test feature selection methods.

Attribute Information

Input variables (based on physicochemical tests): 1 - fixed acidity 2 - volatile acidity 3 - citric acid 4 - residual sugar 5 - chlorides 6 - free sulfur dioxide 7 - total sulfur dioxide 8 - density 9 - pH 10 - sulphates 11 - alcohol Output variable (based on sensory data): 12 - quality (score between 0 and 10)

What might be an interesting thing to do, is to set an arbitrary cutoff for your dependent variable (wine quality) at e.g. 7 or higher getting classified as 'good/1' and the remainder as 'not good/0'. This allows you to practice with hyper parameter tuning on e.g. decision tree algorithms looking at the ROC curve and the AUC value.

You need to build a classification model.

Inspiration

Use machine learning to determine which physicochemical properties make a wine 'good'!

Dataset Link-

<https://github.com/dsrsclentist/DSData/blob/master/winequality-red.csv>

```
In [41]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
```

```
In [42]: url_link = "https://raw.githubusercontent.com/dsrsclentist/DSData/master/winequality-red.csv"
data = pd.read_csv(url_link)
print(data.head(5))
print(data.tail(5))

  fixed acidity  volatile acidity  citric acid  residual sugar  chlorides \
0           7.4             0.70         0.00             1.9         0.076
1           7.8             0.88         0.00             2.6         0.098
2           7.8             0.76         0.04             2.3         0.092
3          11.2             0.28         0.56             1.9         0.075
4           7.4             0.70         0.00             1.9         0.076

  free sulfur dioxide  total sulfur dioxide  density  pH  sulphates \
0             11.0             34.0  0.9978  3.51         0.56
1             25.0             67.0  0.9968  3.20         0.68
2             15.0             54.0  0.9970  3.26         0.65
3             17.0             60.0  0.9980  3.16         0.58
4             11.0             34.0  0.9978  3.51         0.56

  alcohol  quality
0       9.4         5
1       9.8         5
2       9.8         5
3       9.8         6
4       9.4         5

  fixed acidity  volatile acidity  citric acid  residual sugar  chlorides \
1594           6.2             0.600         0.08             2.0         0.090
1595           5.9             0.550         0.10             2.2         0.062
1596           6.3             0.510         0.13             2.3         0.076
1597           5.9             0.645         0.12             2.0         0.075
1598           6.0             0.310         0.47             3.6         0.067

  free sulfur dioxide  total sulfur dioxide  density  pH  sulphates \
1594             32.0             44.0  0.99490  3.45         0.58
1595             39.0             51.0  0.99512  3.52         0.76
1596             29.0             40.0  0.99574  3.42         0.75
1597             32.0             44.0  0.99547  3.57         0.71
1598             18.0             42.0  0.99549  3.39         0.66

  alcohol  quality
1594     10.5         5
1595     11.2         6
1596     11.0         6
1597     10.2         5
1598     11.0         6
```

```
In [43]: print(data.describe())

      fixed acidity  volatile acidity  citric acid  residual sugar \
count  1599.000000      1599.000000  1599.000000  1599.000000
mean    8.319637      0.527821    0.270976    2.538806
std     1.741096      0.179060    0.194801    1.409928
min     4.600000      0.120000    0.000000    0.900000
25%     7.100000      0.390000    0.090000    1.900000
50%     7.900000      0.520000    0.260000    2.200000
75%     9.200000      0.640000    0.420000    2.600000
max    15.900000     1.580000    1.000000   15.500000

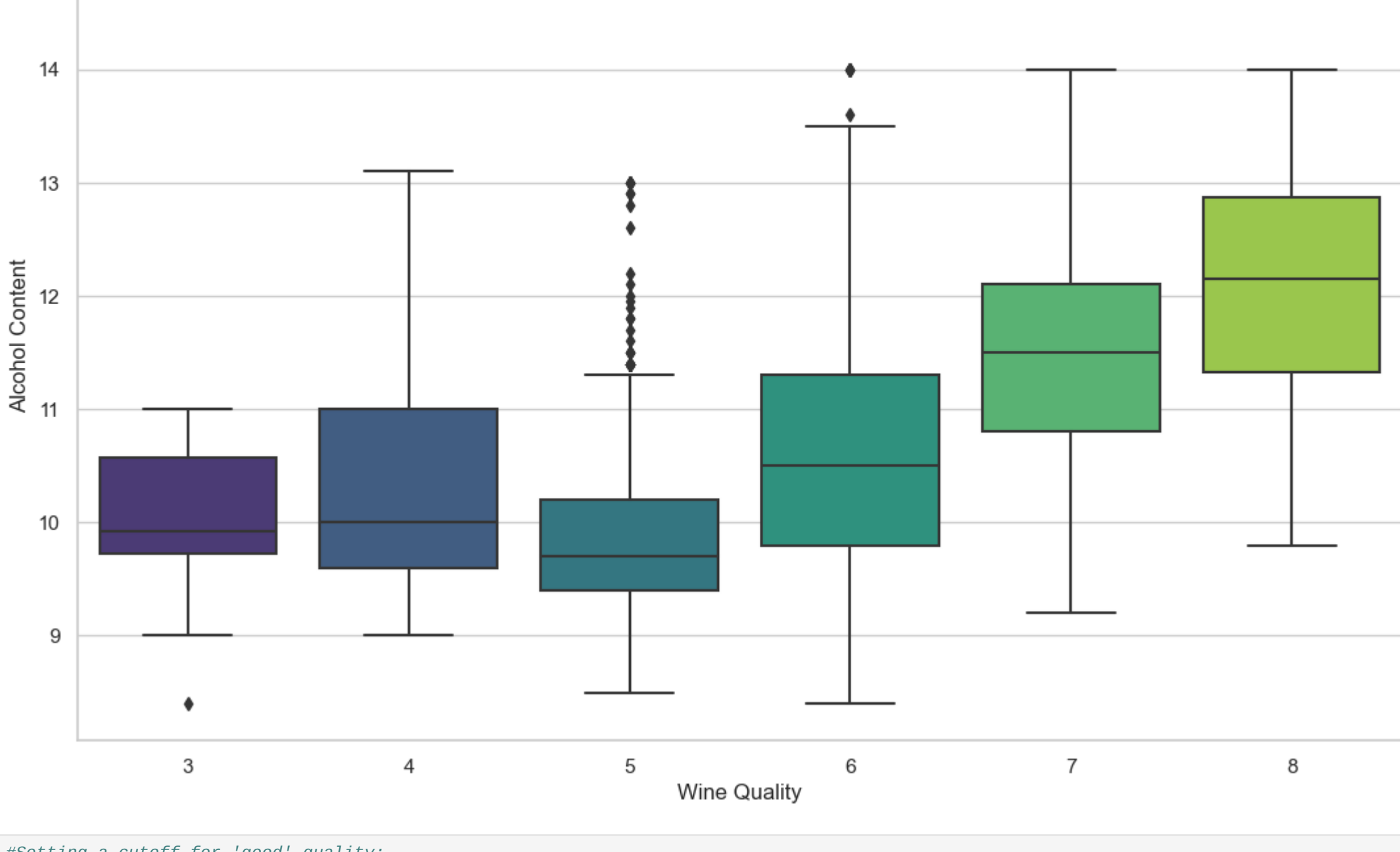
      chlorides  free sulfur dioxide  total sulfur dioxide  density \
count  1599.000000      1599.000000  1599.000000  1599.000000
mean    0.087467     15.874922    46.467792    0.996747
std     0.047065     10.460157    32.895324    0.001887
min     0.012000      1.000000      6.000000    0.990070
25%     0.070000      7.000000     22.000000    0.995600
50%     0.079000     14.000000     38.000000    0.996750
75%     0.090000     21.000000     62.000000    0.997835
max     0.611000     72.000000    289.000000    1.003690

      pH  sulphates  alcohol  quality
count  1599.000000  1599.000000  1599.000000  1599.000000
mean    3.311113    0.658149    10.422983    5.636023
std     0.154386    0.169507    1.065668    0.807569
min     2.740000    0.330000     8.400000    3.000000
25%     3.210000    0.550000     9.500000    5.000000
50%     3.310000    0.620000    10.200000    6.000000
75%     3.400000    0.730000    11.100000    6.000000
max     4.010000    2.000000    14.900000    8.000000
```

```
In [44]: print(data.info())

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1599 entries, 0 to 1598
Data columns (total 12 columns):
#   Column              Non-Null Count  Dtype
---  -
0   fixed acidity        1599 non-null   float64
1   volatile acidity     1599 non-null   float64
2   citric acid          1599 non-null   float64
3   residual sugar       1599 non-null   float64
4   chlorides            1599 non-null   float64
5   free sulfur dioxide  1599 non-null   float64
6   total sulfur dioxide 1599 non-null   float64
7   density              1599 non-null   float64
8   pH                  1599 non-null   float64
9   sulphates            1599 non-null   float64
10  alcohol              1599 non-null   float64
11  quality              1599 non-null   int64
dtypes: float64(11), int64(1)
memory usage: 150.0 KB
None
```

```
In [45]: plt.figure(figsize=(13, 8))
sns.boxplot(x='quality', y='alcohol', data=data, palette='viridis')
plt.title('Distribution of Alcohol Content by Wine Quality')
plt.xlabel('Wine Quality')
plt.ylabel('Alcohol Content')
plt.show()
```



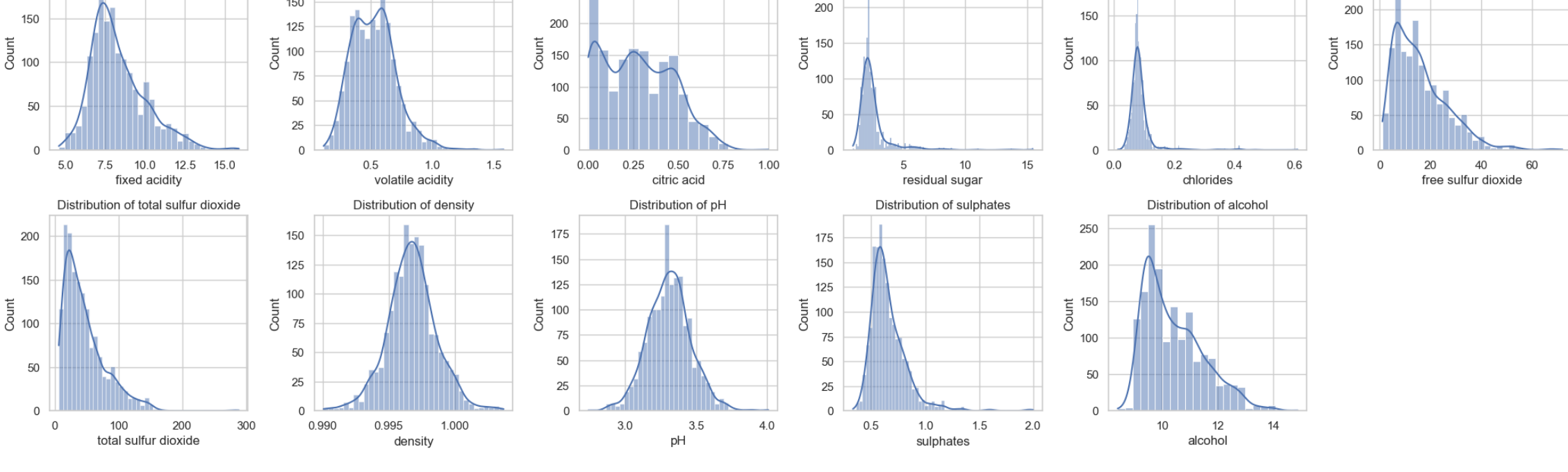
```
In [46]: #Setting a cutoff for 'good' quality:

cutoff = 5
data['quality'] = data['quality'].apply(lambda x: 1 if x >= cutoff else 0)
```

```
In [47]: # Separate features (X) and target variable (y)
X = data.drop('quality', axis=1)
y = data['quality']
```

```
In [48]: # Split 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)
```

```
In [49]: # Visualize the distribution of each input variable
plt.figure(figsize=(21, 10))
for i, column in enumerate(X.columns, 1):
    plt.subplot(3, 6, i)
    sns.histplot(data[column], kde=True)
    plt.title(f'Distribution of {column}')
plt.tight_layout()
plt.show()
```



```
In [50]: # Initialize the Random Forest Classifier
model = RandomForestClassifier(n_estimators=100, random_state=42)
```

```
In [51]: # Train the model
model.fit(X_train, y_train)
```

```
Out[51]: RandomForestClassifier(random_state=42)
```

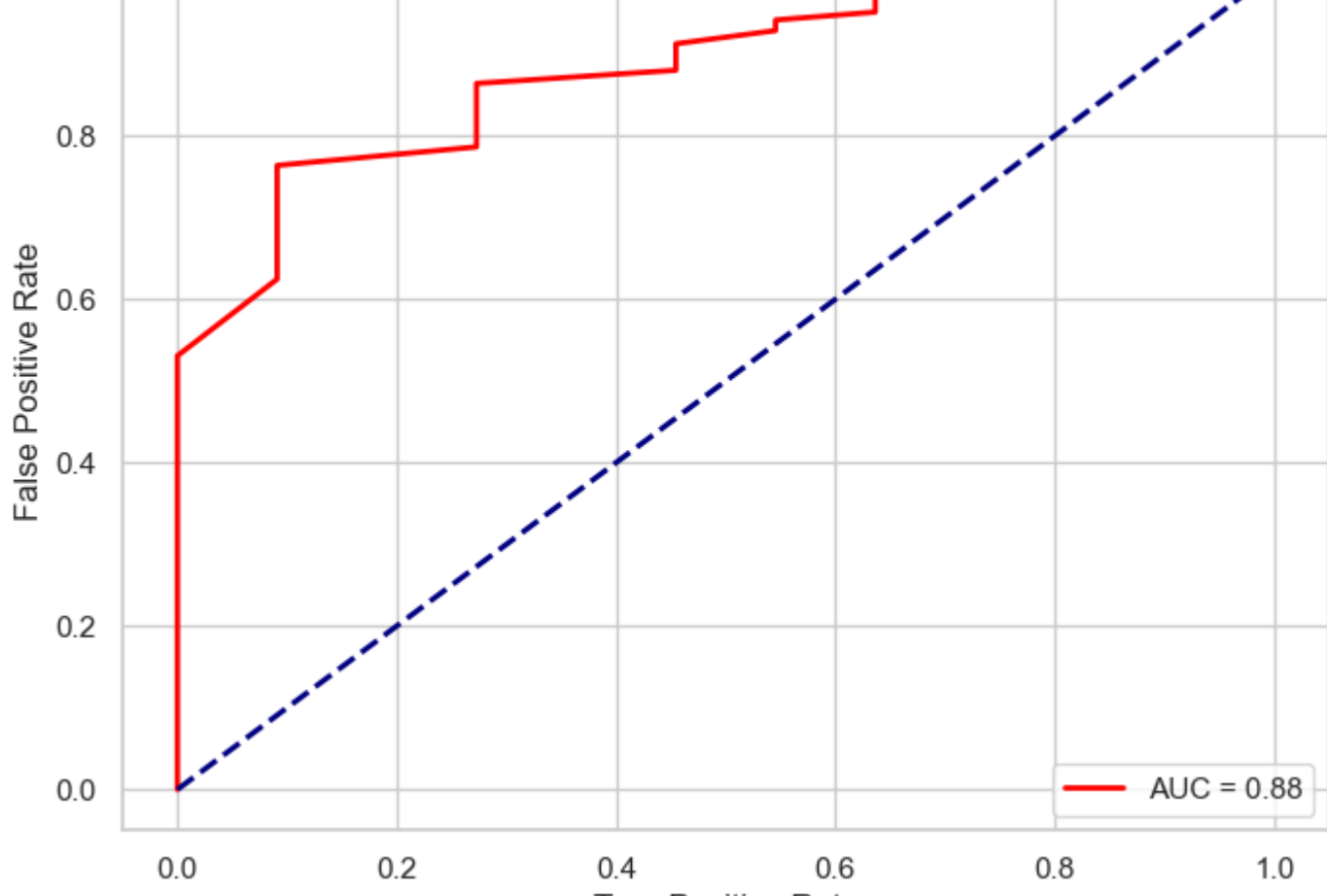
```
In [52]: # Make predictions on the test set
y_pred = model.predict(X_test)

# Evaluate accuracy
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy:.2f}')
```

Accuracy: 0.97

```
In [53]: # Plot ROC curve
plt.figure(figsize=(8, 6))
fpr, tpr, thresholds = roc_curve(y_test, model.predict_proba(X_test)[:, 1])
roc_auc = auc(fpr, tpr)

plt.plot(fpr, tpr, color='red', lw=2, label=f'AUC = {roc_auc:.2f}')
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlabel('True Positive Rate')
plt.ylabel('False Positive Rate')
plt.title(' (ROC) Curve')
plt.legend(loc='lower right')
plt.show()
```



```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```