

```
In [1]: %matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import seaborn as sns

from sklearn.model_selection import train_test_split
from sklearn import tree
from sklearn.naive_bayes import GaussianNB
from sklearn import metrics
```

```
In [2]: pip install seaborn --upgrade
```

```
Requirement already satisfied: seaborn in d:\app\anaconda\lib\site-packages (0.13.2)
Requirement already satisfied: numpy!=1.24.0,>=1.20 in d:\app\anaconda\lib\site-packages (from seaborn) (1.24.3)
Requirement already satisfied: pandas>=1.2 in d:\app\anaconda\lib\site-packages (from seaborn) (2.1.4)
Requirement already satisfied: matplotlib!=3.6.1,>=3.4 in d:\app\anaconda\lib\site-packages (from seaborn) (3.8.0)
Requirement already satisfied: contourpy>=1.0.1 in d:\app\anaconda\lib\site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (1.2.0)
Requirement already satisfied: cycler>=0.10 in d:\app\anaconda\lib\site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (0.11.0)
Requirement already satisfied: fonttools>=4.22.0 in d:\app\anaconda\lib\site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (4.25.0)
Requirement already satisfied: kiwisolver>=1.0.1 in d:\app\anaconda\lib\site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (1.4.4)
Requirement already satisfied: packaging>=20.0 in d:\app\anaconda\lib\site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (23.1)
Requirement already satisfied: pillow>=6.2.0 in d:\app\anaconda\lib\site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (10.2.0)
Requirement already satisfied: pyparsing>=2.3.1 in d:\app\anaconda\lib\site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (3.0.9)
Requirement already satisfied: python-dateutil>=2.7 in d:\app\anaconda\lib\site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in d:\app\anaconda\lib\site-packages (from pandas>=1.2->seaborn) (2023.3.post1)
Requirement already satisfied: tzdata>=2022.1 in d:\app\anaconda\lib\site-packages (from pandas>=1.2->seaborn) (2023.3)
Requirement already satisfied: six>=1.5 in d:\app\anaconda\lib\site-packages (from python-dateutil>=2.7->matplotlib!=3.6.1,>=3.4->seaborn) (1.16.0)
Note: you may need to restart the kernel to use updated packages.
```

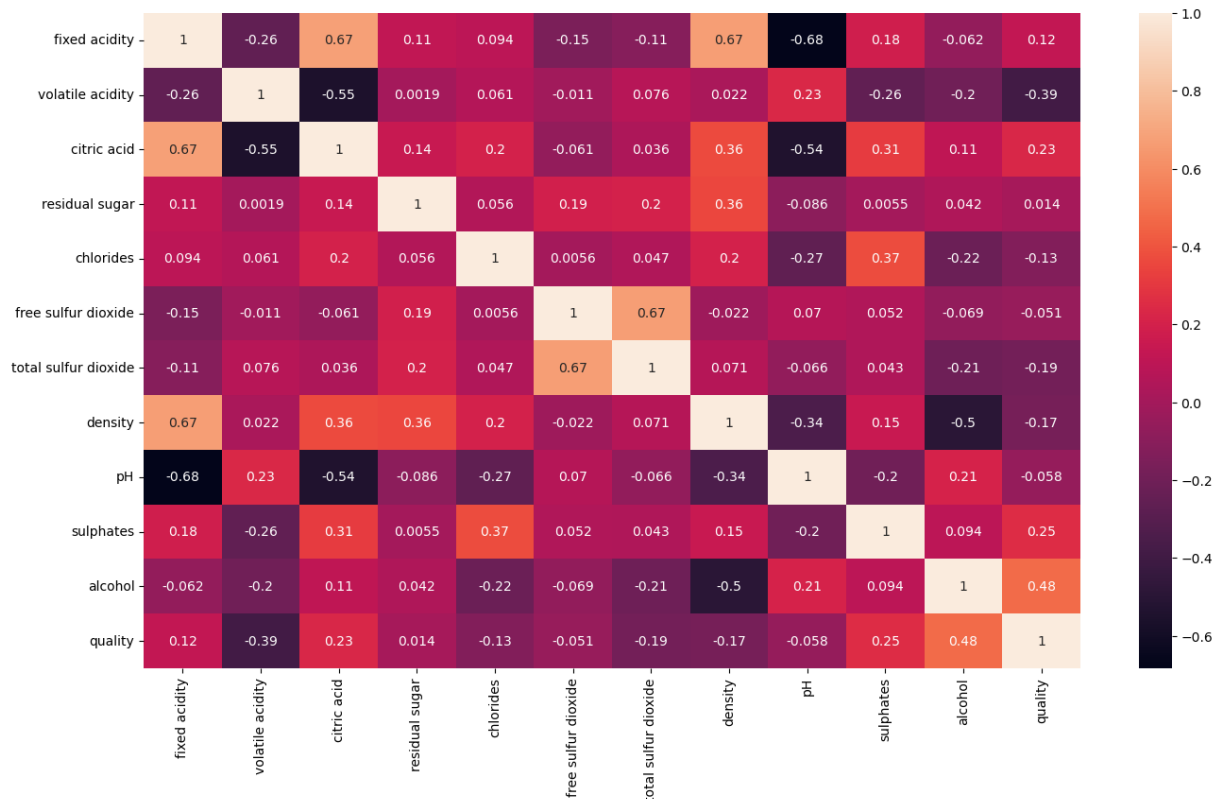
```
In [3]: df = pd.read_csv('winequality-red.csv')
```

```
In [4]: df.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
```

```
In [5]: plt.figure(figsize=(16,9))
sns.heatmap(df.corr(method='pearson'),annot=True)
```

```
Out[5]: <Axes: >
```



```
In [6]: feature = df.drop('quality',axis=1)
label = df['quality']
```

```
In [7]: feature.select_dtypes(exclude=['int64']).columns
```

```
Out[7]: Index(['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar',
              'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density',
              'pH', 'sulphates', 'alcohol'],
              dtype='object')
```

```
In [8]: feature_onehot = pd.get_dummies(feature, columns=feature.select_dtypes(exclude=['int',
feature_onehot
```

```
Out[8]:
```

	fixed acidity_4.6	fixed acidity_4.7	fixed acidity_4.9	fixed acidity_5.0	fixed acidity_5.1	fixed acidity_5.2	fixed acidity_5.3	a
0	False	False	False	False	False	False	False	
1	False	False	False	False	False	False	False	
2	False	False	False	False	False	False	False	
3	False	False	False	False	False	False	False	
4	False	False	False	False	False	False	False	
...	...	...	...	...	...	...	...	
1594	False	False	False	False	False	False	False	
1595	False	False	False	False	False	False	False	
1596	False	False	False	False	False	False	False	
1597	False	False	False	False	False	False	False	
1598	False	False	False	False	False	False	False	

1599 rows × 1453 columns

```
In [9]: x_train, x_test, y_train, y_test = train_test_split(feature, label, test_size=0.3, random
```

```
In [10]: clf = tree.DecisionTreeClassifier(criterion="entropy", random_state=0)
clf.fit(x_train, y_train)
```

```
Out[10]: ▼ DecisionTreeClassifier
DecisionTreeClassifier(criterion='entropy', random_state=0)
```

```
In [11]: tree_pred = clf.predict(x_test)

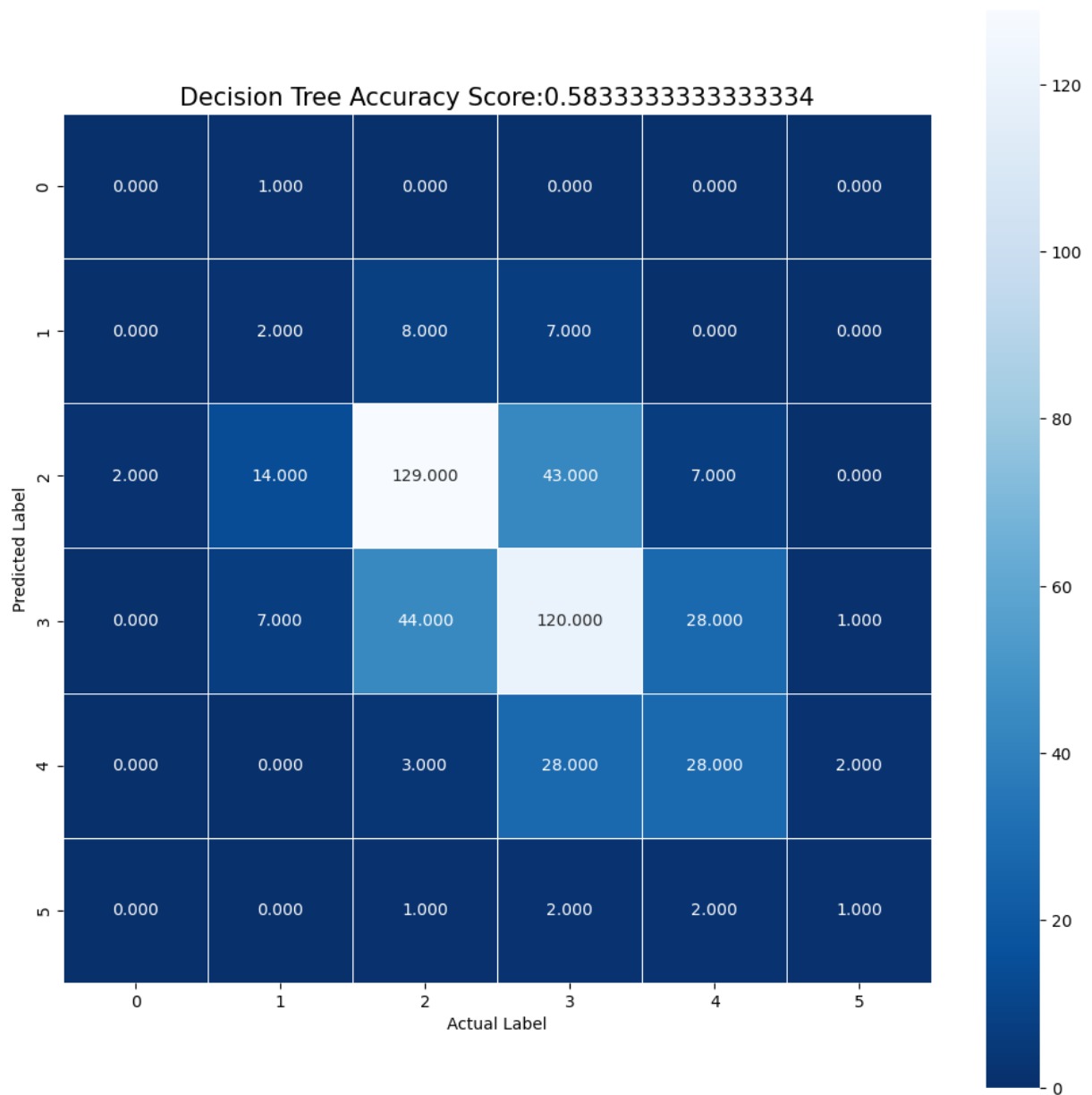
tree_score = metrics.accuracy_score(y_test, tree_pred)
print("Accuracy:", tree_score)
print("Report:", metrics.classification_report(y_test, tree_pred))
```

Accuracy: 0.5833333333333334

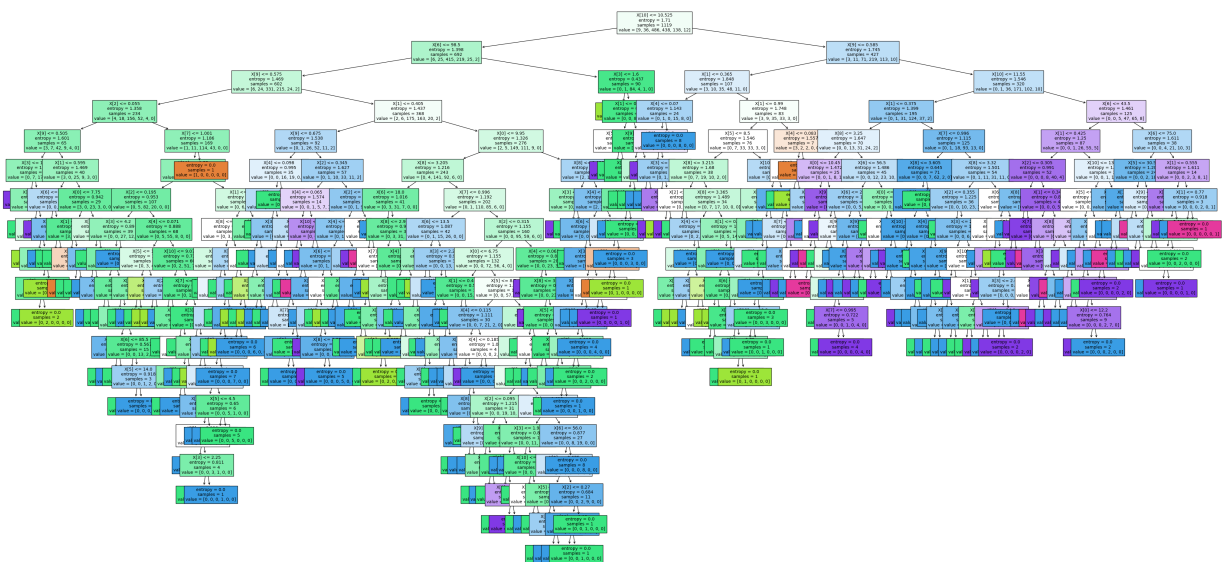
Report:	precision	recall	f1-score	support
3	0.00	0.00	0.00	1
4	0.08	0.12	0.10	17
5	0.70	0.66	0.68	195
6	0.60	0.60	0.60	200
7	0.43	0.46	0.44	61
8	0.25	0.17	0.20	6
accuracy			0.58	480
macro avg	0.34	0.33	0.34	480
weighted avg	0.59	0.58	0.59	480

```
In [12]: tree_cm = metrics.confusion_matrix(y_test, tree_pred)
```

```
In [13]: plt.figure(figsize=(12,12))
sns.heatmap(tree_cm,annot=True, fmt=".3f",linewidth=.5,square=True,cmap='Blues_r');
plt.xlabel('Actual Label');
plt.ylabel('Predicted Label');
title = 'Decision Tree Accuracy Score:{0}'.format(tree_score)
plt.title(title,size=15);
```



```
In [14]: fig, ax = plt.subplots(figsize=(50,24))
tree.plot_tree(clf,filled=True,fontsize=10)
plt.savefig('decision_tree',dpi=100)
plt.show()
```



```
In [15]: clf = tree.DecisionTreeClassifier(criterion="gini",random_state=0)
         clf.fit(x_train,y_train)
```

```
Out[15]: ▼      DecisionTreeClassifier
         DecisionTreeClassifier(random_state=0)
```

```
In [16]: tree_pred = clf.predict(x_test)
         tree_score = metrics.accuracy_score(y_test,tree_pred)
         print("Accuracy:",tree_score)
         print("Report:",metrics.classification_report(y_test,tree_pred))
```

Accuracy: 0.5625

Report:                      precision      recall      f1-score      support

3	0.00	0.00	0.00	1
4	0.07	0.06	0.06	17
5	0.64	0.62	0.63	195
6	0.57	0.58	0.57	200
7	0.45	0.49	0.47	61
8	0.33	0.33	0.33	6

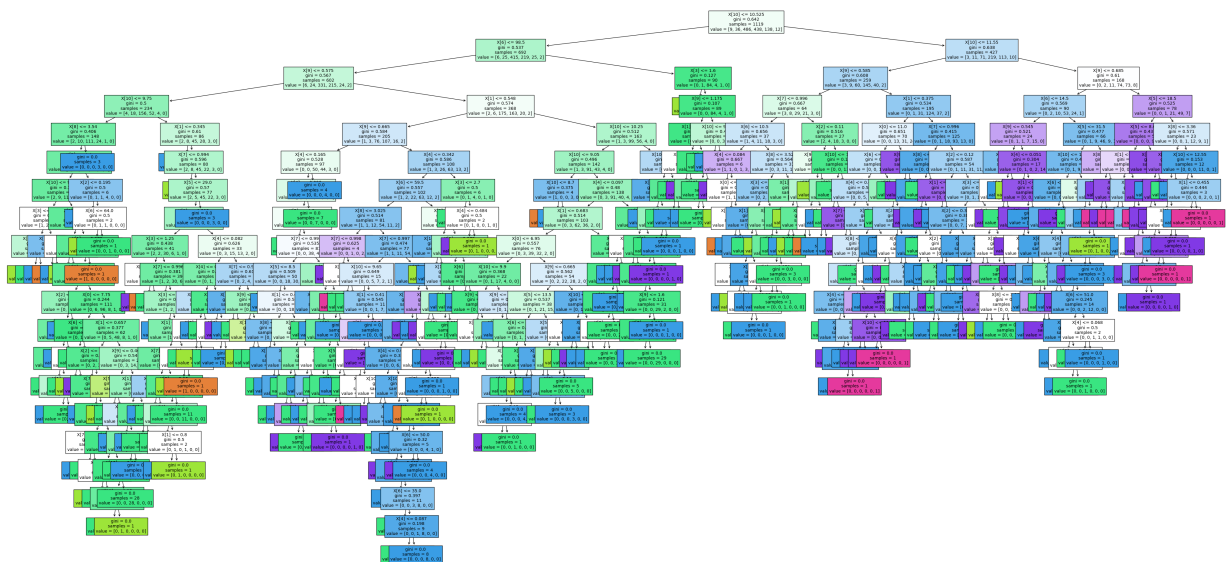
accuracy			0.56	480
macro avg	0.34	0.35	0.35	480
weighted avg	0.56	0.56	0.56	480

```
In [17]: tree_cm = metrics.confusion_matrix(y_test,tree_pred)
```

```
In [18]: plt.figure(figsize=(12,12))
         sns.heatmap(tree_cm,annot=True, fmt=".3f",linewidth=.5,square=True,cmap='Blues_r');
         plt.xlabel('Actual Label');
         plt.ylabel('Predicted Label');
         title = 'Decision Tree Accuracy Score:{0}'.format(tree_score)
         plt.title(title,size=15);
```



```
In [19]: fig, ax = plt.subplots(figsize=(50,24))
tree.plot_tree(clf,filled=True,fontsize=10)
plt.savefig('decision_tree',dpi=100)
plt.show()
```



```
In [20]: gnb = GaussianNB()
```

```
In [21]: bayes_pred = gnb.fit(x_train, y_train).predict(x_test)
bayes_score = metrics.accuracy_score(y_test, bayes_pred)
print("Accuracy: ", bayes_score)
print("Report: ", metrics.classification_report(y_test, bayes_pred))
```

Accuracy: 0.5416666666666666

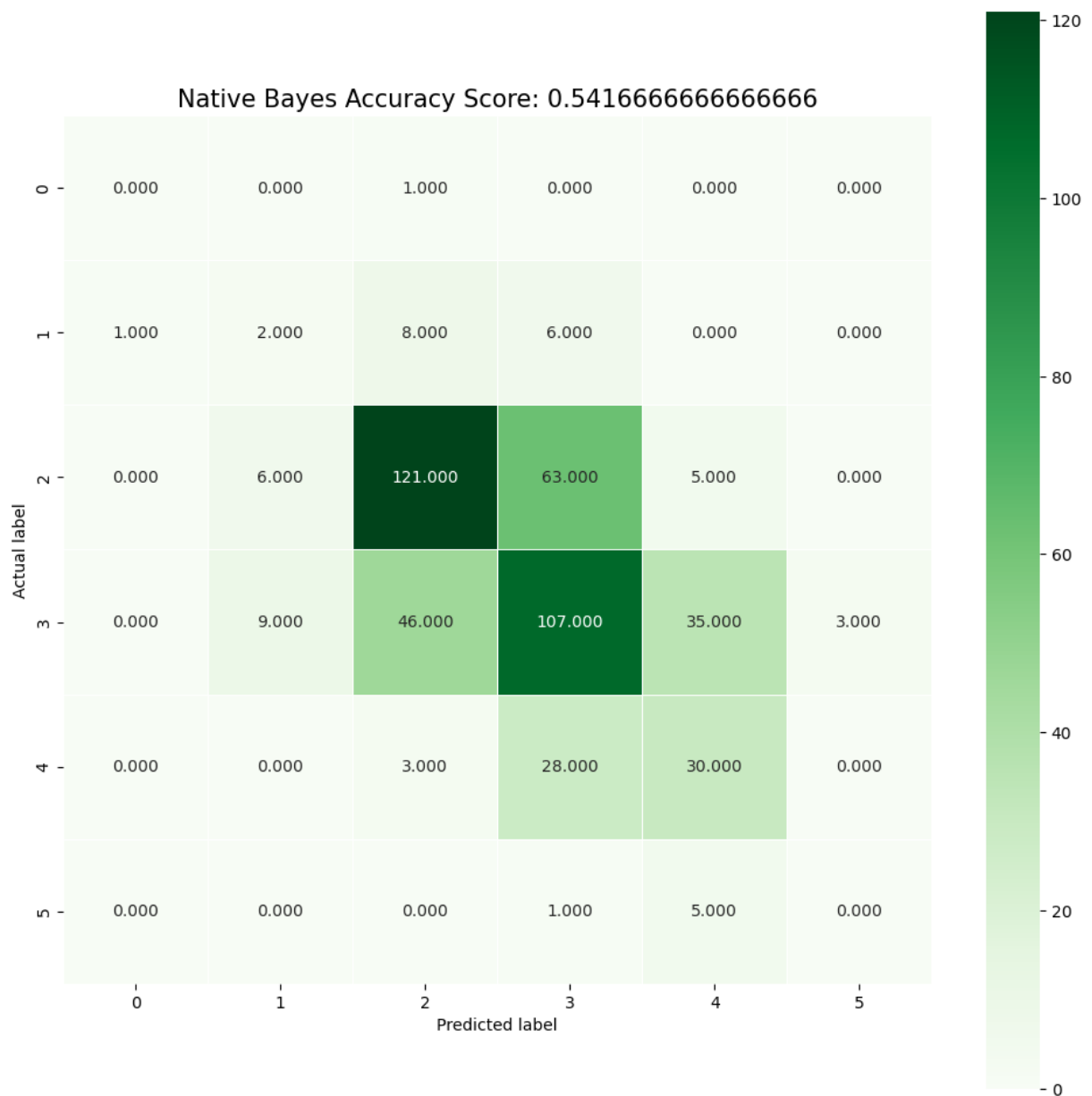
Report: precision recall f1-score support

3	0.00	0.00	0.00	1
4	0.12	0.12	0.12	17
5	0.68	0.62	0.65	195
6	0.52	0.54	0.53	200
7	0.40	0.49	0.44	61
8	0.00	0.00	0.00	6

accuracy			0.54	480
macro avg	0.29	0.29	0.29	480
weighted avg	0.55	0.54	0.54	480

```
In [22]: bayes_cm = metrics.confusion_matrix(y_test, bayes_pred)
plt.figure(figsize=(12,12))
sns.heatmap(bayes_cm,annot=True, fmt=".3f",linewidth=.5,square=True,cmap='Greens');
plt.ylabel('Actual label');
plt.xlabel('Predicted label');
title = 'Native Bayes Accuracy Score: {0}'.format(bayes_score)
plt.title(title, size=15);
```





10 So sánh kết quả của các mô hình trên.

Dựa vào mô hình ta có độ chính xác của các thuật toán lần lượt là:

Thuật toán cây ID3 với 58.34%  
 Thuật toán Naive Bayes với 56.25%  
 Thuật toán cây CART với 54.167%

Vậy đối với mô hình này sử dụng thuật toán cây quyết định ID3 cho ra độ chính xác cao nhất

In [ ]: