

lda-wine-quality

December 19, 2023

```
[52]: # Import necessary libraries
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.metrics import accuracy_score, confusion_matrix, \
    classification_report
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV
```

```
[53]: df = pd.read_csv("/kaggle/input/red-wine-quality-cortez-et-al-2009/
    winequality-red.csv")
df.head()
```

```
[53]:
```

	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

```
[54]: df.columns
```

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

```
[55]: 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
```

```
[56]: df.describe()
```

```
[56]:
```

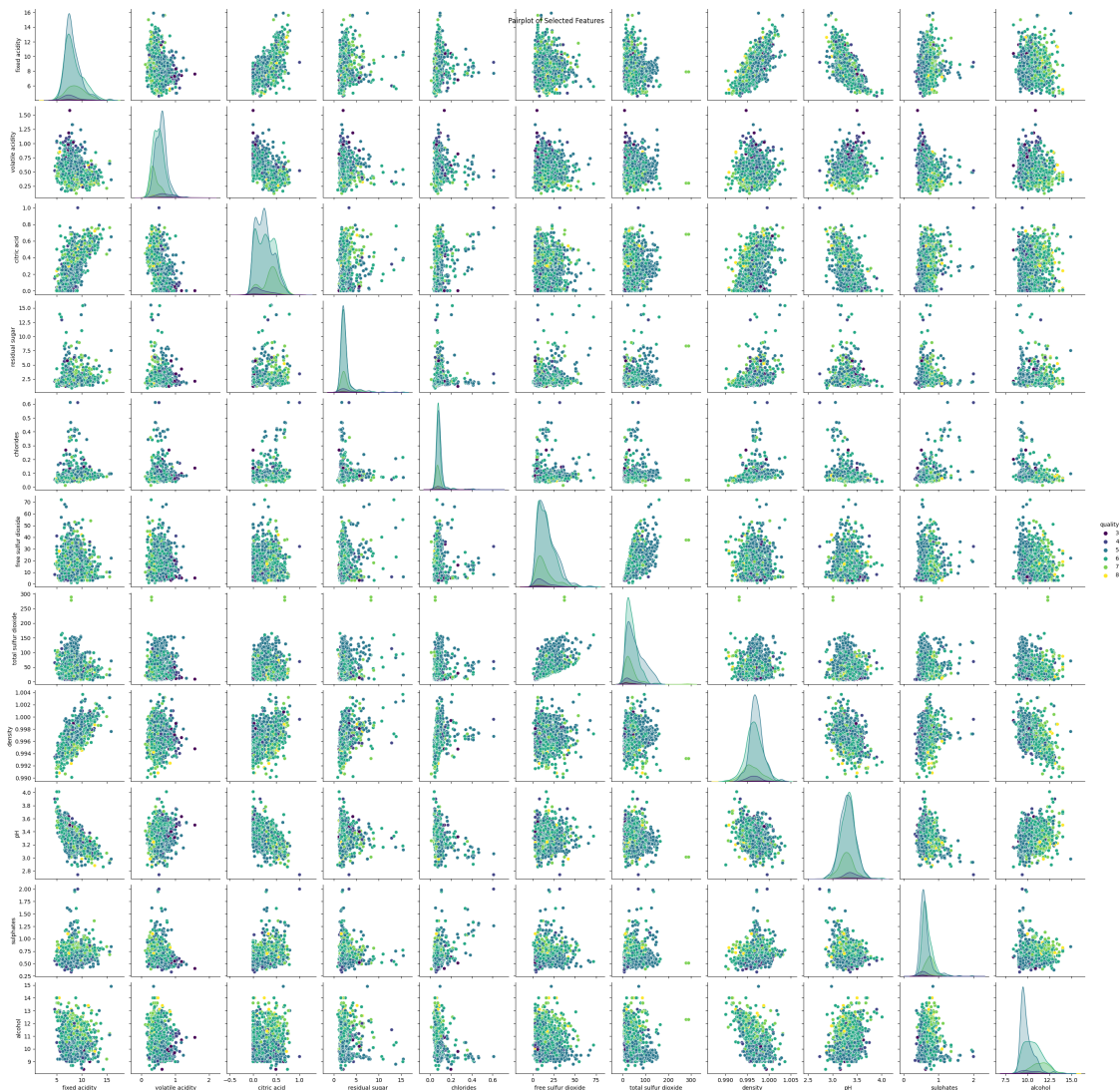
	fixed acidity	volatile acidity	citric acid	residual sugar	\
count	1599.000000	1599.000000	1599.000000	1599.000000	
mean	8.319637	0.527821	0.270976	2.538806	
std	1.741096	0.179060	0.194801	1.409928	
min	4.600000	0.120000	0.000000	0.900000	
25%	7.100000	0.390000	0.090000	1.900000	
50%	7.900000	0.520000	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	

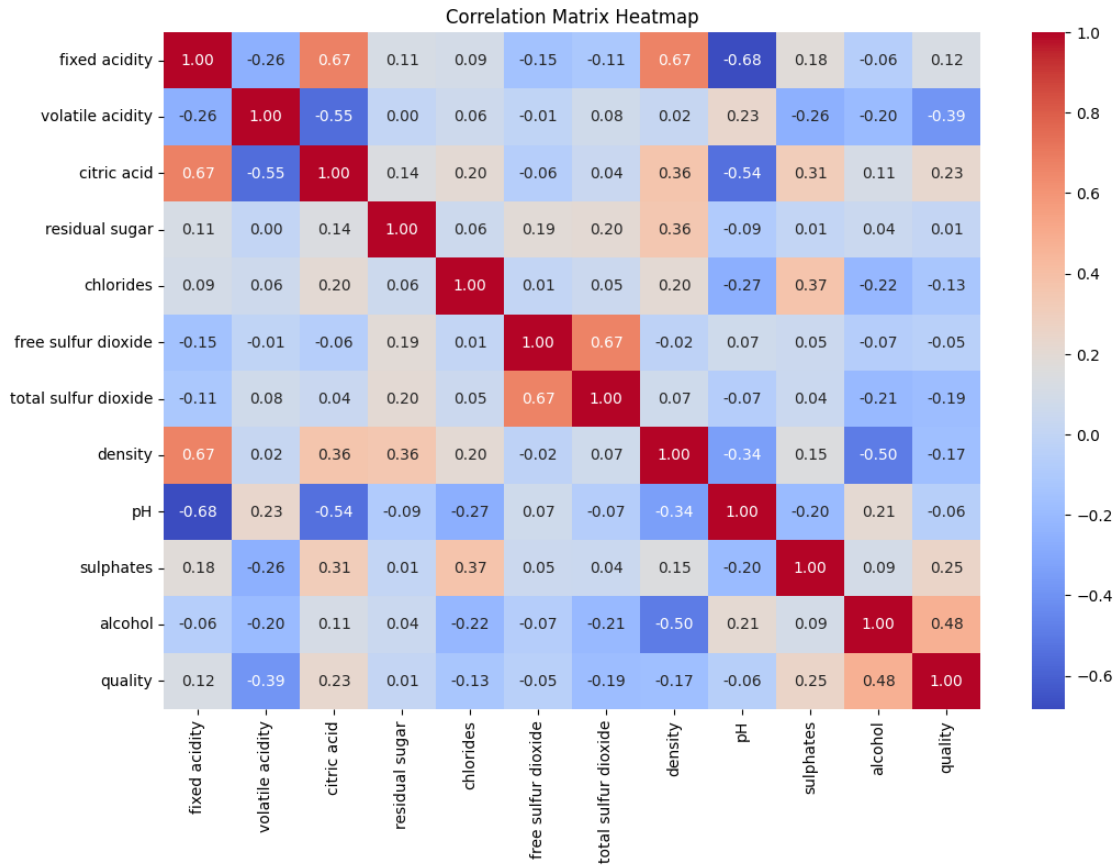
50%	0.079000	14.000000	38.000000	0.996750
75%	0.090000	21.000000	62.000000	0.997835
max	0.611000	72.000000	289.000000	1.003690

	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

```
[57]: # Pairplot for selected features
features = ['fixed acidity', 'volatile acidity', 'citric acid', 'residual_
↪sugar', 'chlorides', 'free sulfur dioxide',
           'total sulfur dioxide', 'density', 'pH', 'sulphates', 'alcohol',_
↪'quality']
sns.pairplot(df[features], hue='quality', markers='o', palette='viridis')
plt.suptitle("Pairplot of Selected Features")
plt.show()
```



```
[58]: # Correlation matrix heatmap
corr_matrix = df.corr()
plt.figure(figsize=(12, 8))
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt=".2f")
plt.title("Correlation Matrix Heatmap")
plt.show()
```



```
[59]: # Data preprocessing
X = df.drop('quality', axis=1)
y = df['quality']

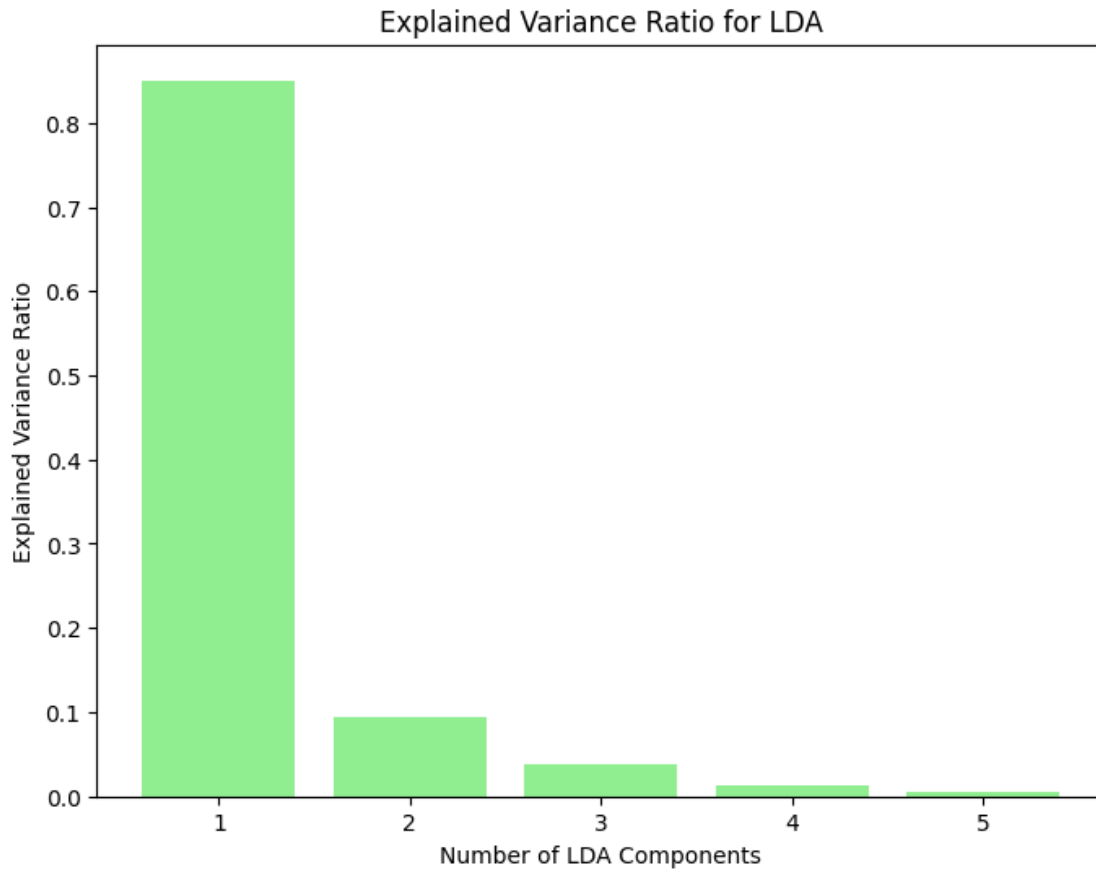
[60]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
↳ random_state=42)

[61]: # Standardize the features
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

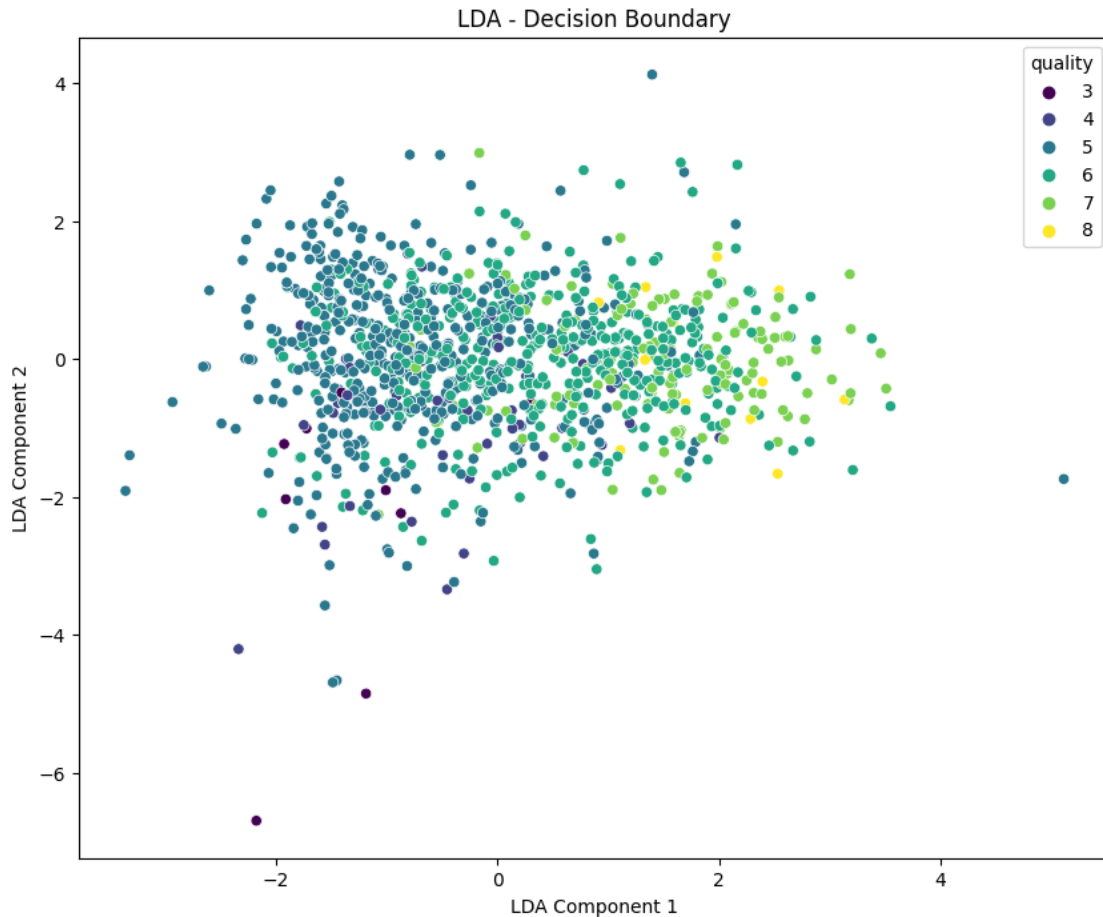
[62]: # Linear Discriminant Analysis (LDA)
lda = LinearDiscriminantAnalysis()
X_train_lda = lda.fit_transform(X_train_scaled, y_train)
X_test_lda = lda.transform(X_test_scaled)

[63]: # Plot the explained variance ratio
plt.figure(figsize=(8, 6))
```

```
plt.bar(range(1, lda.explained_variance_ratio_.shape[0] + 1), lda.
        explained_variance_ratio_, color="lightgreen")
plt.xlabel('Number of LDA Components')
plt.ylabel('Explained Variance Ratio')
plt.title('Explained Variance Ratio for LDA')
plt.show()
```



```
[64]: # Plot the decision boundary
plt.figure(figsize=(10, 8))
sns.scatterplot(x=X_train_lda[:, 0], y=X_train_lda[:, 1], hue=y_train,
                palette='viridis')
plt.title('LDA - Decision Boundary')
plt.xlabel('LDA Component 1')
plt.ylabel('LDA Component 2')
plt.show()
```



```
[65]: # Random Forest Classifier
rf_classifier = RandomForestClassifier(random_state=42)

# Hyperparameter tuning using GridSearchCV
param_grid = {
    'n_estimators': [50, 100, 150],
    'max_depth': [None, 10, 20, 30],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4]
}
```

```
[66]: grid_search = GridSearchCV(estimator=rf_classifier, param_grid=param_grid,
    ↪cv=5, scoring='accuracy')
grid_search.fit(X_train_lda, y_train)
```

```
[66]: GridSearchCV(cv=5, estimator=RandomForestClassifier(random_state=42),
    param_grid={'max_depth': [None, 10, 20, 30],
    'min_samples_leaf': [1, 2, 4],
```

```

        'min_samples_split': [2, 5, 10],
        'n_estimators': [50, 100, 150]},
        scoring='accuracy')

```

```

[67]: # Best parameters from the grid search
best_params = grid_search.best_params_

# Use the best model for prediction
best_rf_model = grid_search.best_estimator_
y_pred = best_rf_model.predict(lda.transform(X_test_scaled))

```

```

[68]: # Model evaluation
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy:.2f}")

```

Accuracy: 0.64

```

[69]: # Classification Report
print("Classification Report:")
print(classification_report(y_test, y_pred))

```

Classification Report:

	precision	recall	f1-score	support
3	0.00	0.00	0.00	1
4	0.00	0.00	0.00	10
5	0.70	0.77	0.74	130
6	0.62	0.64	0.63	132
7	0.52	0.52	0.52	42
8	0.00	0.00	0.00	5
accuracy			0.64	320
macro avg	0.31	0.32	0.31	320
weighted avg	0.61	0.64	0.63	320

```

/opt/conda/lib/python3.10/site-packages/sklearn/metrics/_classification.py:1344:
UndefinedMetricWarning: Precision and F-score are ill-defined and being set to
0.0 in labels with no predicted samples. Use `zero_division` parameter to
control this behavior.

```

```

_warn_prf(average, modifier, msg_start, len(result))

```

```

/opt/conda/lib/python3.10/site-packages/sklearn/metrics/_classification.py:1344:
UndefinedMetricWarning: Precision and F-score are ill-defined and being set to
0.0 in labels with no predicted samples. Use `zero_division` parameter to
control this behavior.

```

```

_warn_prf(average, modifier, msg_start, len(result))

```

```

/opt/conda/lib/python3.10/site-packages/sklearn/metrics/_classification.py:1344:
UndefinedMetricWarning: Precision and F-score are ill-defined and being set to

```


0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

```
_warn_prf(average, modifier, msg_start, len(result))
```

```
[70]: # Confusion Matrix
conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(conf_matrix)
```

Confusion Matrix:

```
[[ 0  0  1  0  0  0]
 [ 0  0  7  3  0  0]
 [ 0  0 100 28  2  0]
 [ 0  0 33 84 15  0]
 [ 0  0  1 19 22  0]
 [ 0  0  0  2  3  0]]
```

```
[71]: plt.figure(figsize=(8, 6))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='bwr',
            xticklabels=sorted(y_test.unique()),
            yticklabels=sorted(y_test.unique()))
plt.title('Confusion Matrix')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
```

