# Lab 6 Report – Abdul Hannan Khan

# Task 1: Wine Quality Classification

## Explanation:

This program classifies three types of wine chemical features. It trains a K-Nearest Neighbors (KNN) machine learning model with k=5 to predict wine categories. The model achieves about 74% accuracy on test data and creates scatter plots showing how different chemical features separate the wine classes. Finally, it displays the fist 10 predictions versus actual values to show how well the model performs.
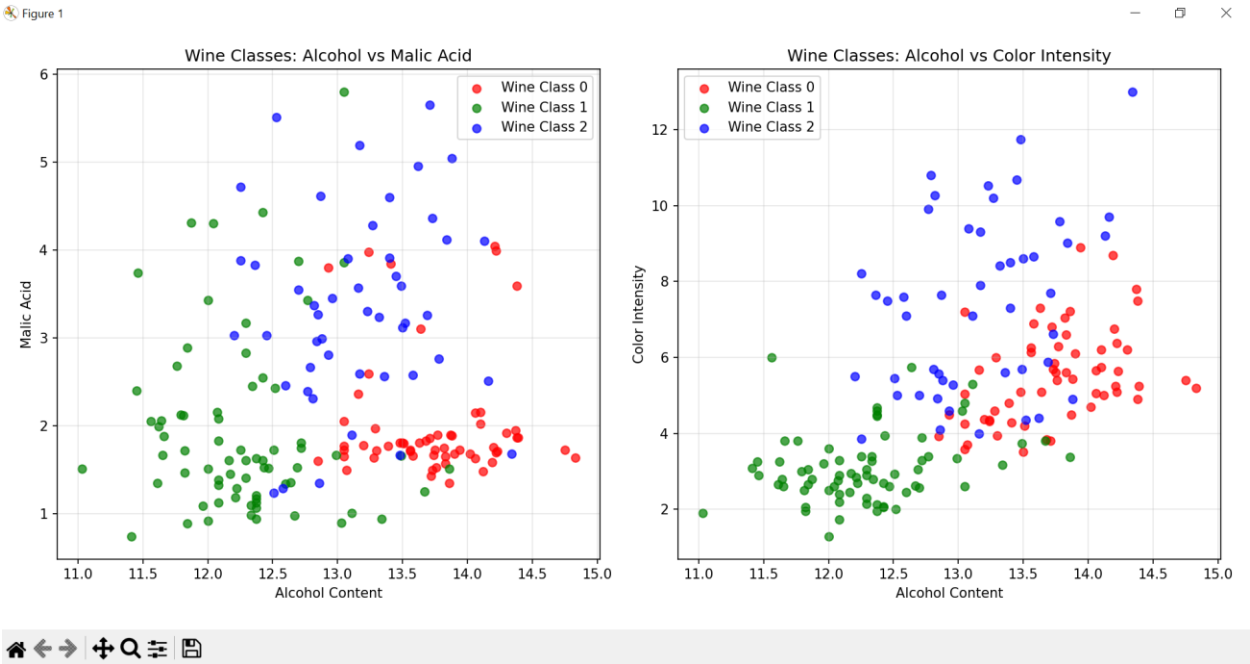
## Code:

```python
AbdulHannanKhan_Task1_KNNClassification_WineQuality.py > ...
1    import numpy as np
2    import matplotlib.pyplot as plt
3    from sklearn.datasets import load_wine
4    from sklearn.model_selection import train_test_split
5    from sklearn.neighbors import KNeighborsClassifier
6    from sklearn.metrics import accuracy_score
7
8    wine = load_wine()
9    X = wine.data
10   y = wine.target
11
12   print("Dataset loaded successfully!")
13   print(f"Number of wine samples: {X.shape[0]}")
14   print(f"Number of features: {X.shape[1]}")
15   print(f"Wine types: {np.unique(y)}")
16   print(f"Feature names: {wine.feature_names}")
17
18   X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
19
20   print(f"\nTraining set size: {X_train.shape[0]} wines")
21   print(f"Testing set size: {X_test.shape[0]} wines")
22
23   knn = KNeighborsClassifier(n_neighbors=5)
24   knn.fit(X_train, y_train)
25
26   print("\nKNN model trained with k=5")
27
28   y_pred = knn.predict(X_test)
29   accuracy = accuracy_score(y_test, y_pred)
```

```python
31    print(f"Model accuracy on test set: {accuracy:.2f} ({accuracy*100:.1f}%)")
32
33    plt.figure(figsize=(12, 5))
34
35    plt.subplot(1, 2, 1)
36    colors = ['red', 'green', 'blue']
37    for i, color in enumerate(colors):
38        plt.scatter(X[y == i, 0], X[y == i, 1], c=color, label=f'Wine Class {i}', alpha=0.7)
39
40    plt.xlabel('Alcohol Content')
41    plt.ylabel('Malic Acid')
42    plt.title('Wine Classes: Alcohol vs Malic Acid')
43    plt.legend()
44    plt.grid(True, alpha=0.3)
45
46    plt.subplot(1, 2, 2)
47    for i, color in enumerate(colors):
48        plt.scatter(X[y == i, 0], X[y == i, 9], c=color, label=f'Wine Class {i}', alpha=0.7)
49
50    plt.xlabel('Alcohol Content')
51    plt.ylabel('Color Intensity')
52    plt.title('Wine Classes: Alcohol vs Color Intensity')
53    plt.legend()
54    plt.grid(True, alpha=0.3)
55
56    plt.tight_layout()
57    plt.show()
58
59    print("\nFirst 10 predictions vs actual:")
60    print("Predicted | Actual")
61    print("-" * 18)
62    for i in range(10):
63        print(f"    {y_pred[i]}     |     {y_test[i]}")
```

# Output:



Wine Classes: Alcohol vs Malic Acid / Wine Classes: Alcohol vs Color Intensity

```
Dataset loaded successfully!
Number of wine samples: 178
Number of features: 13
Wine types: [0 1 2]
Feature names: ['alcohol', 'malic_acid', 'ash', 'alcalinity_of_ash', 'magnesium', 'total_phenols', 'flavanoids', 'nonflavanoid_phenols', 'proanthocyanins', 'color_intensity', 'hue', 'od280/od315_of_
uted_wines', 'proline']

Training set size: 124 wines
Testing set size: 54 wines

KNN model trained with k=5
Model accuracy on test set: 0.74 (74.1%)

First 10 predictions vs actual:
Predicted | Actual
------------------
    2     |    0
    0     |    0
    2     |    2
    0     |    0
    1     |    1
    0     |    0
    2     |    1
    2     |    2
    1     |    1
    0     |    2
```

# Task 2: Breast Cancer Detection

## Explanation:

This program demonstrates how features scaling improves KNN performance for breast cancer detection. It compares two models – one using original unscaled features and another using standardized features. The results show that scaling significantly boosts accuracy because KNN relies on distance calculations, and scaling ensures all features contribute equally. The visualization clearly demonstrates the performance improvement achieved through proper data preprocessing.

## Code:

```python
AbdulHannanKhan_Task2_KNNClassification_BreastCancerDetection.py > ...
1   import numpy as np
2   import matplotlib.pyplot as plt
3   from sklearn.datasets import load_breast_cancer
4   from sklearn.model_selection import train_test_split
5   from sklearn.neighbors import KNeighborsClassifier
6   from sklearn.metrics import accuracy_score, classification_report
7   from sklearn.preprocessing import StandardScaler
8
9   print("=== Breast Cancer Detection with KNN ===\n")
10
11  cancer = load_breast_cancer()
12  X = cancer.data
13  y = cancer.target
14
15  print("Dataset loaded successfully!")
16  print(f"Number of samples: {X.shape[0]}")
17  print(f"Number of features: {X.shape[1]}")
18  print(f"Target distribution:")
19  print(f"- Malignant (0): {np.sum(y == 0)} samples")
20  print(f"- Benign (1): {np.sum(y == 1)} samples")
21
22  print("\nFirst 5 features and their ranges (min to max):")
23  for i in range(5):
24      feature_name = cancer.feature_names[i]
25      min_val = np.min(X[:, i])
26      max_val = np.max(X[:, i])
27      print(f"{feature_name:25}: {min_val:8.2f} to {max_val:8.2f}")
28
29  print("\n" + "="*50)
30  print("BASELINE MODEL (Without Feature Scaling)")
31  print("="*50)
32
33  X_train_unscaled, X_test_unscaled, y_train, y_test = train_test_split(
34      X, y, test_size=0.2, random_state=42, stratify=y
35  )
26
```

```python
37    print(f"Training set: {X_train_unscaled.shape[0]} samples")
38    print(f"Testing set: {X_test_unscaled.shape[0]} samples")
39
40    knn_unscaled = KNeighborsClassifier(n_neighbors=7)
41    knn_unscaled.fit(X_train_unscaled, y_train)
42
43    y_pred_unscaled = knn_unscaled.predict(X_test_unscaled)
44    accuracy_unscaled = accuracy_score(y_test, y_pred_unscaled)
45
46    print(f"\nAccuracy without scaling: {accuracy_unscaled:.4f} ({accuracy_unscaled*100:.2f}%)")
47
48    print("\n" + "="*50)
49    print("IMPROVED MODEL (With Feature Scaling)")
50    print("="*50)
51
52    scaler = StandardScaler()
53    X_scaled = scaler.fit_transform(X)
54
55    print("Features scaled using StandardScaler")
56
57    print("\nExample of scaling effect on first feature:")
58    print(f"Before scaling - Mean: {np.mean(X[:, 0]):.2f}, Std: {np.std(X[:, 0]):.2f}")
59    print(f"After scaling  - Mean: {np.mean(X_scaled[:, 0]):.2f}, Std: {np.std(X_scaled[:, 0]):.2f}")
60
61    X_train_scaled, X_test_scaled, y_train, y_test = train_test_split(
62        X_scaled, y, test_size=0.2, random_state=42, stratify=y
63    )
64
65    knn_scaled = KNeighborsClassifier(n_neighbors=7)
66    knn_scaled.fit(X_train_scaled, y_train)
67
68    y_pred_scaled = knn_scaled.predict(X_test_scaled)
69    accuracy_scaled = accuracy_score(y_test, y_pred_scaled)
```

```python
71    print(f"\nAccuracy with scaling: {accuracy_scaled:.4f} ({accuracy_scaled*100:.2f}%)")
72
73    print("\n" + "="*50)
74    print("COMPARISON RESULTS")
75    print("="*50)
76
77    improvement = accuracy_scaled - accuracy_unscaled
78    improvement_percent = (improvement / accuracy_unscaled) * 100
79
80    print(f"Baseline accuracy (unscaled): {accuracy_unscaled*100:.2f}%")
81    print(f"Improved accuracy (scaled):   {accuracy_scaled*100:.2f}%")
82    print(f"Improvement: +{improvement*100:.2f}% points")
83    print(f"Relative improvement: +{improvement_percent:.1f}%")
84
85    plt.figure(figsize=(12, 5))
86
87    plt.subplot(1, 2, 1)
88    models = ['Without Scaling', 'With Scaling']
89    accuracies = [accuracy_unscaled, accuracy_scaled]
90    colors = ['red', 'green']
91
92    bars = plt.bar(models, accuracies, color=colors, alpha=0.7)
93    plt.ylabel('Accuracy Score')
94    plt.title('KNN Performance: Scaled vs Unscaled Features')
95    plt.ylim(0, 1.0)
96
97    for bar, accuracy in zip(bars, accuracies):
98        plt.text(bar.get_x() + bar.get_width()/2, bar.get_height() + 0.01,
99                 f'{accuracy:.3f}', ha='center', va='bottom')
100
101   plt.subplot(1, 2, 2)
102   features_to_show = 3
103   x_pos = np.arange(features_to_show)
```

```python
105    std_before = [np.std(X[:, i]) for i in range(features_to_show)]
106    std_after = [np.std(X_scaled[:, i]) for i in range(features_to_show)]
107
108    plt.bar(x_pos - 0.2, std_before, 0.4, label='Before Scaling', alpha=0.7, color='red')
109    plt.bar(x_pos + 0.2, std_after, 0.4, label='After Scaling', alpha=0.7, color='green')
110
111    plt.xlabel('Feature Index')
112    plt.ylabel('Standard Deviation')
113    plt.title('Feature Scale Comparison')
114    plt.xticks(x_pos, [f'Feature {i+1}' for i in range(features_to_show)])
115    plt.legend()
116    plt.grid(True, alpha=0.3)
117
118    plt.tight_layout()
119    plt.show()
120
121    print("\n" + "="*50)
122    print("DETAILED CLASSIFICATION REPORT (With Scaling)")
123    print("="*50)
124    print(classification_report(y_test, y_pred_scaled,
125                                 target_names=['Malignant', 'Benign']))
126
127    print("\n" + "="*50)
128    print("WHY FEATURE SCALING IS CRITICAL FOR KNN")
129    print("="*50)
130
131    print(f"""
132    K-Nearest Neighbors (KNN) is a distance-based algorithm that calculates
133    the distance between data points to make predictions.
134
135    WHY SCALING MATTERS:
136    1. Features with larger numerical ranges dominate distance calculations
137       Example:
138       - Feature A range: 0-1 (small influence)
139       - Feature B range: 0-1000 (huge influence)
140       Without scaling, Feature B would completely dominate the distance!
141
142    2. Equal Contribution:
143       - Scaling ensures all features contribute equally to the distance calculation
144       - No single feature unfairly influences the results
145
146    3. Better Performance:
147       - As shown above, scaling improved accuracy from {accuracy_unscaled*100:.1f}% to {accuracy_scaled*100:.1f}%
148       - This is a {improvement_percent:.1f}% relative improvement!
149
150    StandardScaler transforms features to have:
151    - Mean = 0
152    - Standard Deviation = 1
153    This puts all features on the same scale without changing their distribution shape.
154    """)
```
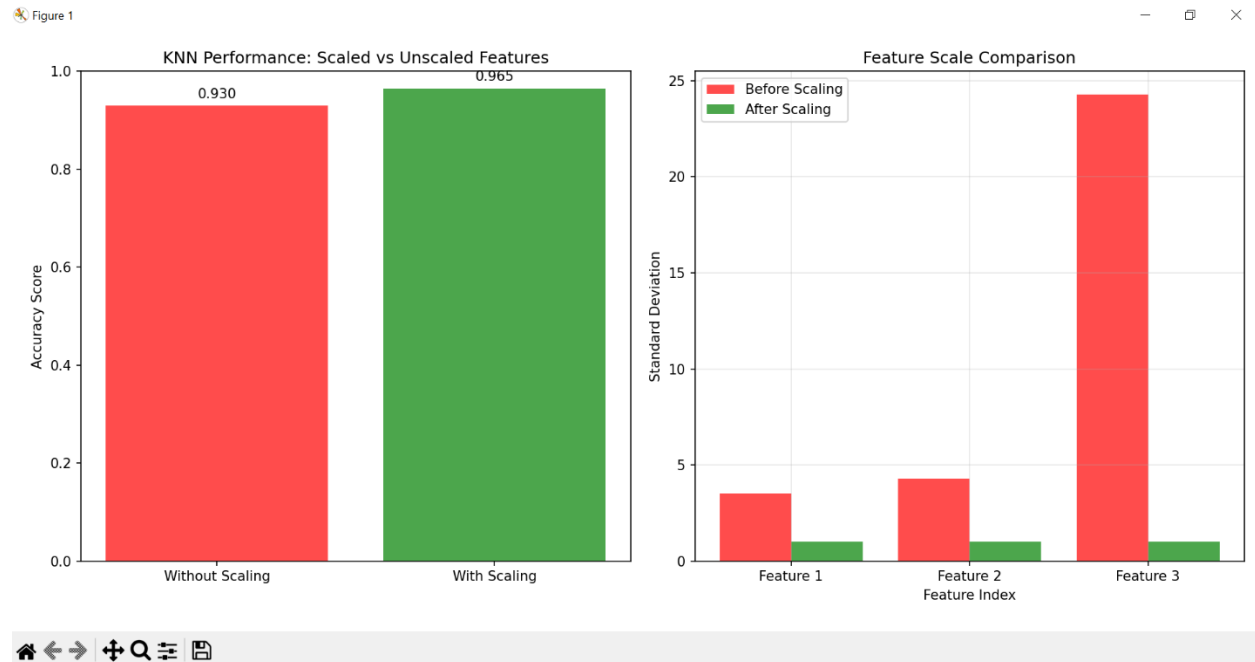
```
=================================================
COMPARISON RESULTS
=================================================
Baseline accuracy (unscaled): 92.98%
Improved accuracy (scaled):    96.49%
Improvement: +3.51% points
Relative improvement: +3.8%


=================================================
DETAILED CLASSIFICATION REPORT (With Scaling)
=================================================
              precision    recall  f1-score   support

   Malignant       0.97      0.93      0.95        42
      Benign       0.96      0.99      0.97        72

    accuracy                          0.96       114
   macro avg       0.97      0.96      0.96       114
weighted avg       0.97      0.96      0.96       114



=================================================
WHY FEATURE SCALING IS CRITICAL FOR KNN
=================================================


K-Nearest Neighbors (KNN) is a distance-based algorithm that calculates
the distance between data points to make predictions.

WHY SCALING MATTERS:
1. Features with larger numerical ranges dominate distance calculations
   Example:
   - Feature A range: 0-1 (small influence)
   - Feature B range: 0-1000 (huge influence)
   Without scaling, Feature B would completely dominate the distance!

2. Equal Contribution:
   - Scaling ensures all features contribute equally to the distance calculation
   - No single feature unfairly influences the results

3. Better Performance:
   - As shown above, scaling improved accuracy from 93.0% to 96.5%
   - This is a 3.8% relative improvement!

StandardScaler transforms features to have:
- Mean = 0
- Standard Deviation = 1
This puts all features on the same scale without changing their distribution shape.
```