

Assignment 8-9:

Total Points: 2 points)

Objective:

Use a real-world dataset to analyze and compare Decision Trees, Random Forest, Ensemble Models (Bagging and Boosting), Dimensionality Reduction (PCA), and Clustering (k-means). This assignment aims to explore the effectiveness of these methods, understand their strengths and limitations, and identify practical use cases.

Dataset:

Use the "Heart Disease Dataset" from Kaggle: <https://www.kaggle.com/datasets/fedesoriano/heart-failure-prediction>. This dataset contains health indicators such as age, cholesterol levels, and smoking habits, with the goal of predicting the likelihood of heart disease.

Instructions:

1. Dataset Preparation (0.5 points):

Download and explore the Heart Disease Dataset.

Preprocess the data:

Handle missing values (if any).

Normalize/standardize the numerical features.

Encode categorical variables.

Visualize key variables to understand relationships (e.g., age vs. cholesterol).

2. Task 1: Decision Trees and Random Forest (0.5 points):

Define a classification task, such as predicting whether a person has heart disease ("Yes" or "No").

Implement Decision Tree and Random Forest models:

Perform hyperparameter tuning for both models.

Record performance metrics such as accuracy, precision, recall, and F1-score.

Extract and visualize feature importance from both models.

Write a short analysis comparing Decision Tree and Random Forest, focusing on overfitting, interpretability, and predictive performance.

3. Task 2: Ensemble Models – Bagging and Boosting (0.5 points):

Implement Bagging and Boosting methods using:

Bagging: BaggingClassifier (e.g., with Decision Trees).

Boosting: AdaBoost or Gradient Boosting.

Evaluate and compare the performance of these models against Random Forest.

Discuss how Bagging and Boosting handle bias-variance trade-offs.

4. Task 3: Dimensionality Reduction and Clustering (0.5 points):

Dimensionality Reduction:

Apply PCA to the dataset and retain components that explain 90% of the variance.

Visualize the dataset in 2D using the top two principal components.

Clustering:

Perform k-means clustering on the dataset (or PCA-reduced data).

Determine the optimal number of clusters using the elbow method or silhouette analysis.

Compare the clusters to the actual heart disease labels and evaluate clustering performance using metrics like purity or Adjusted Rand Index.

Grading (2 Points Total):

Dataset Preparation and Visualization (0.5 points):

Thorough preprocessing and meaningful visualizations of relationships between features.

Implementation and Results (1 point):

Accurate implementation of all models and dimensionality reduction techniques.

Clear evaluation of results with metrics and visualizations.

Analysis and Comparison (0.5 points):

Insightful comparison of methods, highlighting trade-offs and practical recommendations.

```
In [13]: import pandas as pd

In [4]: df = pd.read_csv('C:\\Users\\Acer\\Downloads\\heart.csv')

In [5]: (df.head())

Out[6]:
   Age  Sex  ChestPainType  RestingBP  Cholesterol  FastingBS  RestingECG  MaxHR  ExerciseAngina  Oldpeak  ST_Slope  HeartDisease
0  40   M      ATA         140         289           0      Normal      172             N      0.0      Up           0
1  49   F      NAP         160         180           0      Normal      156             N      1.0      Flat           1
2  37   M      ATA         130         283           0      ST         98             N      0.0      Up           0
3  46   F      ASY         138         214           0      Normal      108             Y      1.5      Flat           1
4  54   M      NAP         140         195           0      Normal      122             N      0.0      Up           0

In [7]: (df.info())

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 918 entries, 0 to 917
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Age                   918 non-null    int64
 1   Sex                   918 non-null    object
 2   ChestPainType         918 non-null    object
 3   RestingBP             918 non-null    int64
 4   Cholesterol            918 non-null    int64
 5   FastingBS             918 non-null    int64
 6   RestingECG            918 non-null    object
 7   MaxHR                 918 non-null    int64
 8   ExerciseAngina        918 non-null    object
 9   Oldpeak               918 non-null    float64
10   ST_Slope              918 non-null    object
11   HeartDisease          918 non-null    int64
dtypes: float64(1), int64(6), object(5)
memory usage: 16.4 MB

In [8]: (df.isnull().sum())

Out[8]:
Age                0
Sex                0
ChestPainType      0
RestingBP          0
Cholesterol         0
FastingBS          0
RestingECG         0
MaxHR              0
ExerciseAngina     0
Oldpeak            0
ST_Slope           0
HeartDisease       0
dtype: int64

In [12]: from sklearn.preprocessing import StandardScaler

# List of numerical columns to standardize
numerical_columns = ['Age', 'RestingBP', 'Cholesterol', 'FastingBS', 'MaxHR', 'Oldpeak']

# Initialize StandardScaler
scaler = StandardScaler()

# Standardize the numerical features
df[numerical_columns] = scaler.fit_transform(df[numerical_columns])

# Check the standardized data
print(df[numerical_columns].head())

In [13]: from sklearn.preprocessing import LabelEncoder

# Initialize LabelEncoder
label_encoder = LabelEncoder()

# Encode categorical features
df['Sex'] = label_encoder.fit_transform(df['Sex'])
df['ChestPainType'] = label_encoder.fit_transform(df['ChestPainType'])
df['RestingECG'] = label_encoder.fit_transform(df['RestingECG'])
df['ExerciseAngina'] = label_encoder.fit_transform(df['ExerciseAngina'])
df['ST_Slope'] = label_encoder.fit_transform(df['ST_Slope'])

# Check the encoded data
print(df.head())

In [15]: import matplotlib.pyplot as plt
import seaborn as sns

# Visualize relationships between key features
# Example: Age vs Cholesterol
plt.figure(figsize=(8,6))
sns.scatterplot(x=Age, y=Cholesterol, data=df, hue=HeartDisease)
plt.title('Age vs Cholesterol for Heart Disease Prediction')
plt.show()

# Plot histograms for numerical features
df[numerical_columns].hist(figsize=(10, 8), bins=30)
plt.suptitle('Distribution of Numerical Features')
plt.show()
```



