

# Assignment-4: Decision Tree & XGBooster

Mahbub Ahmed Turza  
ID: 2211063042  
North South University  
mahbub.turza@northsouth.edu

November 15, 2024

## Abstract

Machine learning has become a cornerstone of predictive modeling across diverse domains, providing powerful tools for regression, classification, and decision-making tasks. This report investigates the performance of two widely used algorithms—Decision Tree Regressor/Classifier and XGBoost Regressor/Classifier—on four distinct datasets: Car Dataset, DT Brain Cancer, DT-Wage, and DT-Credit. These datasets represent a variety of tasks, including regression, binary classification, and multi-class classification, ensuring a comprehensive evaluation.

The study employs metrics such as Mean Squared Error (MSE), accuracy, and classification reports to assess model performance. Results reveal the trade-offs between simplicity and predictive power, with Decision Trees excelling in interpretability and XGBoost demonstrating superior accuracy and generalization. By analyzing the algorithms' strengths and weaknesses across diverse datasets, this report offers insights into their suitability for different types of predictive tasks. The findings contribute to better-informed decisions when selecting machine learning models for real-world applications.

## I. INTRODUCTION

In recent years, machine learning has become an indispensable tool for predictive modeling in various domains, enabling accurate forecasting and decision-making. This report explores the application of machine learning algorithms on four distinct datasets: Car Dataset, DT Brain Cancer, DT-Wage, and DT-Credit. These datasets encompass a wide range of prediction tasks, from regression and classification to multi-class predictions, making them ideal for evaluating the strengths and weaknesses of different algorithms.

### A. Objective

The primary objective of this study is to compare the performance of two popular machine learning models—**Decision Tree Regressor** and **XGBoost Regressor/Classifier**—across the four datasets. The models are evaluated using key performance metrics such as Mean Squared Error (MSE), R-squared ( $R^2$ ), accuracy, and residual analysis to assess their suitability for each dataset and task. Additionally, the study aims to provide insights into the trade-offs between model interpretability and predictive accuracy.

### B. Datasets

The four datasets used in this study are as follows:

- **DT-Credit:** Multi-class Classification
- **Brain Cancer:** Binary Classification
- **Wage Dataset:** Regression
- **Credit Dataset:** Regression

### *C. Significance of Study*

The selection of these datasets allows for an evaluation of the models under varying conditions and complexities. The diversity in dataset structure, feature types, and target variables provides a robust testbed for understanding how different algorithms handle different types of prediction problems. By analyzing the performance of Decision Trees and XGBoost across these datasets, this study contributes to the growing body of knowledge on the applicability of machine learning models in real-world scenarios.

### *D. Report Structure*

The remainder of this report is structured as follows:

- **Methodology:** Details the preprocessing steps, hyperparameter tuning, and model evaluation strategies.
- **Results:** Presents the performance metrics of the models on each dataset.
- **Discussion:** Compares and contrasts the findings across datasets and models.
- **Conclusion:** Summarizes the key takeaways and provides recommendations for model selection based on task requirements.

This comprehensive approach ensures that the report provides valuable insights into the capabilities and limitations of the chosen machine learning models across diverse datasets and tasks.

## II. METHODOLOGY

### A. For Car Dataset:

- **Data Preprocessing** : To start, we load and preprocess the car dataset, which contains features such as buying price, maintenance cost, number of doors, passenger capacity, trunk size, and safety level. Each feature contributes to determining the car's overall acceptability.

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.preprocessing import LabelEncoder
4
5 # Load dataset
6 data = pd.read_csv('car_data.csv')
7 print(data.head())
8
9 # Initialize label encoder
10 label_encoder = LabelEncoder()
11
12 # Encoding each categorical column
13 for col in data.columns:
14     data[col] = label_encoder.fit_transform(data[col])
15 print(data.head())
16
```

Each categorical feature is converted into numeric form using `LabelEncoder` from `sklearn.preprocessing`, as both models require numerical input.

- **Splitting Data into Training and Testing Sets:** To evaluate model performance, the dataset is split into training and testing sets.

```
1 X = data.drop('class', axis=1)
2 y = data['class']
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
4 random_state=42)
```

Here, the target variable `class` is separated from features `X`. An 80-20 split is used, where 80% of the data is used for training, and 20% is reserved for testing.

- **Model Selection and Training: 1) Decision Tree Classifier:** The Decision Tree classifier is initialized and trained on the training set.

```
1 from sklearn.tree import DecisionTreeClassifier
2
3 dt_model = DecisionTreeClassifier(random_state=42)
4 dt_model.fit(X_train, y_train)
5
```

**2) XGBoost Classifier:** The XGBoost model is initialized and trained with specific configuration parameters.

```
1 from xgboost import XGBClassifier
2
3 xgb_model = XGBClassifier(use_label_encoder=False, eval_metric='mlogloss',
4 random_state=42)
5 xgb_model.fit(X_train, y_train)
```

The `XGBClassifier` is initialized with `use_label_encoder=False` and `eval_metric='mlogloss'`, as required by XGBoost for multi-class classification.

- **Model Evaluation** : The models are evaluated on accuracy, precision, recall, and F1-score.

```
1 from sklearn.metrics import accuracy_score, classification_report
2
3 dt_pred = dt_model.predict(X_test)
```

```

4     xgb_pred = xgb_model.predict(X_test)
5
6     dt_accuracy = accuracy_score(y_test, dt_pred)
7     xgb_accuracy = accuracy_score(y_test, xgb_pred)
8
9     print("Decision Tree Accuracy:", dt_accuracy)
10    print("XGBoost Accuracy:", xgb_accuracy)
11
12    print("Decision Tree Classification Report:")
13    print(classification_report(y_test, dt_pred))
14
15    print("XGBoost Classification Report:")
16    print(classification_report(y_test, xgb_pred))
17

```

accuracy\_score computes the accuracy, while classification\_report provides precision, recall, and F1-score for each class.

- **Comparative Analysis Using Confusion Matrix :** The confusion matrix for each model is displayed for better insight into misclassifications.

```

1     from sklearn.metrics import confusion_matrix
2     import matplotlib.pyplot as plt
3     import seaborn as sns
4
5     dt_conf_matrix = confusion_matrix(y_test, dt_pred)
6     plt.figure(figsize=(8, 6))
7     sns.heatmap(dt_conf_matrix, annot=True, fmt="d", cmap="Blues")
8     plt.title("Decision Tree Confusion Matrix")
9     plt.xlabel("Predicted Label")
10    plt.ylabel("True Label")
11    plt.show()
12
13    xgb_conf_matrix = confusion_matrix(y_test, xgb_pred)
14    plt.figure(figsize=(8, 6))
15    sns.heatmap(xgb_conf_matrix, annot=True, fmt="d", cmap="Greens")
16    plt.title("XGBoost Confusion Matrix")
17    plt.xlabel("Predicted Label")
18    plt.ylabel("True Label")
19    plt.show()
20

```

The confusion matrices are generated using confusion\_matrix from sklearn.metrics, visualized with seaborn and matplotlib. The matrices illustrate true positives, false positives, and false negatives for both classifiers.

## B. For Brain Cancer Dataset

### • Importing Libraries

Essential libraries for data handling, visualization, and model training were imported as follows:

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from sklearn.model_selection import train_test_split, GridSearchCV
6 from sklearn.tree import DecisionTreeClassifier
7 from xgboost import XGBClassifier
8 from sklearn.metrics import precision_recall_curve
9 from sklearn.preprocessing import StandardScaler, OneHotEncoder
10
```

### • Loading and Exploring the Data

The dataset was loaded, checked for missing values, and unnecessary columns were dropped.

```
1 df = pd.read_csv("DT-BrainCancer.csv")
2 df.info() # To understand data types and missing values
3 df = df.dropna() # Removing rows with missing values
4 df = df.drop('Unnamed: 0', axis=1) # Dropping unnecessary column
5
```

### • Encoding Categorical Variables

Categorical variables were encoded using one-hot encoding.

```
1 categorical = df.select_dtypes(include=['object']).columns
2 encoder = OneHotEncoder(drop='first', sparse_output=False)
3 encoded_data = encoder.fit_transform(df[categorical])
4 encoded_df = pd.DataFrame(encoded_data, columns=encoder.get_feature_names_out(
5 categorical))
6 df = pd.concat([df.drop(categorical, axis=1), encoded_df], axis=1)
```

### • Defining Features (X) and Target (y)

The target variable status was separated from the features.

```
1 X = df.drop(columns=['status'])
2 y = df['status']
3
```

### • Scaling the Features

Standardization was applied to scale the features.

```
1 scaler = StandardScaler()
2 X_scaled = scaler.fit_transform(X)
3
```

### • Splitting the Data

Data was split into training, validation, and testing sets.

```
1 X_train, X_temp, y_train, y_temp = train_test_split(X_scaled, y, test_size=0.3,
2 random_state=42)
3 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5,
4 random_state=42)
```

### • Training Decision Tree Model with Hyperparameter Tuning

Grid search was performed over hyperparameters to optimize the Decision Tree model.

```
1 dt = DecisionTreeClassifier(random_state=42)
2 param_grid_dt = {
```

```

3     'max_depth': [5, 10, 15, 20],
4     'min_samples_split': [2, 5, 10],
5     'min_samples_leaf': [1, 2, 4]
6 }
7 grid_search_dt = GridSearchCV(estimator=dt, param_grid=param_grid_dt, cv=5, scoring='
accuracy')
8 grid_search_dt.fit(X_train, y_train)
9 best_dt = grid_search_dt.best_estimator_
10

```

## • Training XGBoost Model with Hyperparameter Tuning

Similar grid search tuning was applied for the XGBoost model.

```

1 xgb_model = XGBClassifier(random_state=42)
2 param_grid_xgb = {
3     'learning_rate': [0.001, 0.01, 0.05, 0.1],
4     'max_depth': [3, 5, 7],
5     'n_estimators': [100, 200, 300]
6 }
7 grid_search_xgb = GridSearchCV(estimator=xgb_model, param_grid=param_grid_xgb, cv=5,
scoring='accuracy')
8 grid_search_xgb.fit(X_train, y_train)
9 best_xgb = grid_search_xgb.best_estimator_
10

```

## • Evaluating Model Performance

We assessed accuracy, confusion matrix, precision, recall, and F1 score for both models.

```

1 def accuracy(y_true, y_pred):
2     return np.sum(y_true == y_pred) / len(y_true)
3
4 def confusion_matrix_manual(y_true, y_pred):
5     tp = np.sum((y_true == 1) & (y_pred == 1))
6     tn = np.sum((y_true == 0) & (y_pred == 0))
7     fp = np.sum((y_true == 0) & (y_pred == 1))
8     fn = np.sum((y_true == 1) & (y_pred == 0))
9     return tp, tn, fp, fn
10

```

## • Comparing Performance Metrics

We made predictions and calculated metrics for both Decision Tree and XGBoost models.

```

1 y_train_pred_dt = best_dt.predict(X_train)
2 y_val_pred_dt = best_dt.predict(X_val)
3 y_test_pred_dt = best_dt.predict(X_test)
4
5 y_train_pred_xgb = best_xgb.predict(X_train)
6 y_val_pred_xgb = best_xgb.predict(X_val)
7 y_test_pred_xgb = best_xgb.predict(X_test)
8

```

## • Precision-Recall Curve

Precision-recall curves were plotted to visualize the trade-off between precision and recall.

```

1 from sklearn.metrics import precision_recall_curve
2
3 def plot_precision_recall_curve(y_true, y_prob, label):
4     precision, recall, _ = precision_recall_curve(y_true, y_prob)
5     plt.plot(recall, precision, label=label)
6     plt.xlabel('Recall')
7     plt.ylabel('Precision')
8     plt.title('Precision-Recall Curve')
9     plt.legend()

```

```
10     plt.show()
11
12     plot_precision_recall_curve(y_test, best_dt.predict_proba(X_test)[:, 1], label='
Decision Tree')
13     plot_precision_recall_curve(y_test, best_xgb.predict_proba(X_test)[:, 1], label='
XGBoost')
```

### C. For Wage Dataset

- **Loading and Exploring the Data :**

```
1 from google.colab import drive
2 drive.mount('/content/drive')
3
4 import pandas as pd
5
6 df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/DT-Wage.csv')
7 df.head()
8 df.shape
9 df.info()
10 df.describe()
11 df.isnull().sum()
12
```

This step involves mounting Google Drive to access the dataset, loading the CSV file using `pd.read_csv()`, and performing initial data exploration using functions like `.head()`, `.shape`, `.info()`, `.describe()`, and `.isnull().sum()` to check for missing values and gain insights into the dataset structure.

- **Analyzing Unique Values in Categorical Columns:**

```
1 un = df['maritl'].unique()
2 print(un)
3 # Similar code for 'race', 'education', 'region', 'jobclass', 'health'
4
```

This part identifies unique values in each categorical column, helping to determine categories in each variable and aiding in the encoding process.

- **Separating Numerical and Categorical Features:**

```
1 categorical = df.select_dtypes(include=['object']).columns
2 numerical = df.select_dtypes(exclude=['object']).columns
3
```

Here, the columns are separated into categorical and numerical variables using `select_dtypes()`, allowing for customized preprocessing for each data type.

- **One-Hot Encoding of Categorical Features:**

```
1 from sklearn.preprocessing import OneHotEncoder
2
3 encoder = OneHotEncoder(drop='first', sparse_output=False)
4 encoded_data = encoder.fit_transform(df[categorical])
5 encoded_df = pd.DataFrame(encoded_data, columns=encoder.get_feature_names_out(
6 categorical))
7 encoded_df.index = df.index
8 df = df.drop(categorical, axis=1)
9 df = pd.concat([df, encoded_df], axis=1)
```

One-hot encoding is used to transform categorical features into binary columns. The transformed features are then merged back with the original DataFrame after dropping the original categorical columns.

- **Separating Features and Target Variable :**

```
1 X = df.drop(columns=['wage'])
2 y = df['wage']
3
```

Here, the feature matrix  $X$  and the target variable  $y$  are separated, with `wage` as the target.

- **Data Standardization:**

```
1 from sklearn.preprocessing import StandardScaler
2
```



```

3     scaler = StandardScaler()
4     X_scaled = scaler.fit_transform(X)
5

```

The features are standardized using `StandardScaler()` to ensure that all variables contribute equally to the model.

- **Train-Validation-Test Split:**

```

1     from sklearn.model_selection import train_test_split
2
3     X_train, X_temp, y_train, y_temp = train_test_split(X_scaled, y, test_size=0.3,
4     random_state=42)
5     X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5,
6     random_state=42)
7

```

The data is split into training, validation, and test sets, with 70% for training, 15% for validation, and 15% for testing.

- **Setting up Parameter Grid and Performing Grid Search:**

```

1     from sklearn.tree import DecisionTreeRegressor
2     from sklearn.model_selection import GridSearchCV
3
4     param_grid = {
5         'max_depth': [3, 5, 7, 10, 12, None],
6         'min_samples_split': [2, 5, 10, 12],
7         'min_samples_leaf': [1, 2, 5, 7],
8         'max_features': ['sqrt', 'log2', None],
9         'criterion': ['squared_error', 'friedman_mse', 'absolute_error'],
10    }
11    model = DecisionTreeRegressor(random_state=42)
12    grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=5, n_jobs=-1,
13    verbose=2, scoring='neg_mean_squared_error')
14    grid_search.fit(X_train, y_train)

```

Hyperparameter tuning is performed using `GridSearchCV` to test different combinations of hyperparameters defined in `param_grid`. This grid search uses 5-fold cross-validation, scoring models based on negative mean squared error.

- **Identifying Best Hyperparameters:**

```

1     best_params = grid_search.best_params_
2     print("Best Parameters:", best_params)
3

```

After grid search, the optimal hyperparameters are identified and displayed.

- **Model Evaluation :**

```

1     from sklearn.metrics import mean_squared_error, r2_score
2
3     best_model = grid_search.best_estimator_
4
5     y_train_pred = best_model.predict(X_train)
6     y_val_pred = best_model.predict(X_val)
7     y_test_pred = best_model.predict(X_test)
8
9     mse_train = mean_squared_error(y_train, y_train_pred)
10    mse_val = mean_squared_error(y_val, y_val_pred)
11    mse_test = mean_squared_error(y_test, y_test_pred)
12
13    r2_train = r2_score(y_train, y_train_pred)
14    r2_val = r2_score(y_val, y_val_pred)

```

```
15     r2_test = r2_score(y_test, y_test_pred)
16
17     print("Mean Squared Error (Train):", mse_train)
18     print("Mean Squared Error (Validation):", mse_val)
19     print("Mean Squared Error (Test):", mse_test)
20     print("R    (Train): {:.2f}".format(r2_train))
21     print("R    (Validation): {:.2f}".format(r2_val))
22     print("R    (Test): {:.2f}".format(r2_test))
23
```

The model's performance is evaluated using Mean Squared Error (MSE) and  $R^2$  scores for the training, validation, and test sets. MSE measures the average squared difference between predicted and actual values, while  $R^2$  reflects the percentage of variance in the target variable explained by the model.

#### D. For Credit Dataset

The objective of this project is to predict the Balance variable of a dataset (DT-Credit.csv), using various machine learning models, particularly focusing on Decision Trees and XGBoost. Below is the step-by-step explanation of the methodology.

- **Data Loading:** The dataset DT-Credit.csv is loaded into a pandas DataFrame using the `pd.read_csv` function. This dataset contains information about individuals, including their financial attributes such as Income, Limit, Rating, etc.

```
1 df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/DT-Credit.csv')
2
```

- **Exploratory Data Analysis (EDA):** The dataset is explored to understand its structure and check for missing values. This is done by examining the shape, data types, and summary statistics of the dataset. We also checked the unique values for categorical columns (Own, Student, Married, and Region) to understand the range of values they can take.

```
1 df.shape
2 df.info()
3 df.describe()
4 df.isnull().sum()
5
```

- **Encoding Categorical Variables :** The categorical features Own, Student, Married, and Region are encoded using one-hot encoding. One-hot encoding transforms categorical variables into numerical format by creating binary columns for each category.

```
1 encoder = OneHotEncoder(drop='first', sparse_output=False)
2 encoded_data = encoder.fit_transform(df[categorical])
3 encoded_df = pd.DataFrame(encoded_data, columns=encoder.get_feature_names_out(
4 categorical))
```

The original categorical columns are dropped from the dataset, and the newly encoded columns are appended to the DataFrame.

```
1 df = df.drop(categorical, axis=1)
2 df = pd.concat([df, encoded_df], axis=1)
3
```

- **Feature and Target Variable Selection:** The features X and the target variable y are defined. The target variable is the Balance, and the features are the rest of the columns in the dataset.

```
1 X = df.drop(columns=['Balance'])
2 y = df['Balance']
3
```

- **Data Preprocessing:** Standard scaling is applied to the feature data X to normalize the values. This ensures that each feature contributes equally to the machine learning models.

```
1 scaler = StandardScaler()
2 X_scaled = scaler.fit_transform(X)
3
```

- **Data Splitting :** The dataset is split into training, validation, and test sets using `train_test_split`. Initially, 70% of the data is used for training, and 30% is split between validation and test sets (15% each).

```
1 X_train, X_temp, y_train, y_temp = train_test_split(X_scaled, y, test_size=0.3,
2 random_state=42)
3 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5,
4 random_state=42)
```

- **Hyperparameter Tuning with GridSearchCV:** We perform hyperparameter tuning to find the best parameters for the models. A grid search is conducted over multiple hyperparameters using GridSearchCV to select the optimal values that minimize the error for Decision Trees.

```
1 param_grid = {
2     'max_depth': [3, 5, 7, 10, 12, None],
3     'min_samples_split': [2, 5, 10, 12],
4     'min_samples_leaf': [1, 2, 5, 7],
5     'max_features': ['sqrt', 'log2', None],
6     'criterion': ['squared_error', 'friedman_mse', 'absolute_error'],
7 }
8
9 grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=5, n_jobs=-1,
10 verbose=2, scoring='neg_mean_squared_error')
11 grid_search.fit(X_train, y_train)
```

- **Model Training:** After determining the best hyperparameters for the DecisionTreeRegressor using GridSearchCV, we fit the model on the training data:

```
1 best_model = grid_search.best_estimator_
2
```

We then make predictions on the training, validation, and test sets.

```
1 y_train_pred = best_model.predict(X_train)
2 y_val_pred = best_model.predict(X_val)
3 y_test_pred = best_model.predict(X_test)
4
```

- **Performance Evaluation:** To evaluate the performance of the model, the Mean Squared Error (MSE) and  $R^2$  score are computed for the training, validation, and test sets. The MSE indicates how well the model is predicting the target variable, with a lower value indicating better performance. The  $R^2$  score indicates how much variance in the target variable is explained by the model, with higher values indicating better model performance.

```
1 mse_train = mean_squared_error(y_train, y_train_pred)
2 mse_val = mean_squared_error(y_val, y_val_pred)
3 mse_test = mean_squared_error(y_test, y_test_pred)
4
5 r2_train = r2_score(y_train, y_train_pred)
6 r2_val = r2_score(y_val, y_val_pred)
7 r2_test = r2_score(y_test, y_test_pred)
8
```

- **XGBoost Model:** After evaluating the Decision Tree model, we used XGBoost (XGBRegressor) to predict the target variable. We applied a similar process, including data splitting, hyperparameter tuning with GridSearchCV, training the model, and evaluating performance using MSE and  $R^2$ .

```
1 xgb_model = XGBRegressor(objective='reg:squarederror', random_state=42)
2
3 grid_search = GridSearchCV(estimator=xgb_model, param_grid=param_grid, cv=5, n_jobs
4 =-1, verbose=2, scoring='neg_mean_squared_error')
5 grid_search.fit(X_train, y_train)
```

We then evaluate the XGBoost model with MSE and  $R^2$ , comparing its performance to the Decision Tree model.

- **Model Comparison and Visualization:** Finally, the predicted values from the models are compared to the true values in scatter plots, showing how well the models predict the target variable. Additionally, residuals are plotted to observe the distribution of errors for each model.

```
1 plt.scatter(y_train, y_train_pred, color="blue", label="Train")
2 plt.scatter(y_val, y_val_pred, color="orange", label="Validation")
3 plt.scatter(y_test, y_test_pred, color="green", label="Test")
4 plt.plot([y_train.min(), y_train.max()], [y_train.min(), y_train.max()], 'k--', lw
5         =2)
```

Residuals are plotted similarly to evaluate the error patterns:

```
1 plt.scatter(y_train_pred, train_residuals, color="blue", label="Train")
2 plt.scatter(y_val_pred, val_residuals, color="orange", label="Validation")
3 plt.scatter(y_test_pred, test_residuals, color="green", label="Test")
4 plt.axhline(0, color="black", linestyle="--", linewidth=2)
5
```

### III. EXPERIMENTS RESULTS

#### A. For Car Dataset

1) *Dataset Description:* The experiment was conducted using the Car Evaluation dataset, which contains six categorical attributes related to car buying preferences: buying price, maintenance cost, number of doors, passenger capacity, luggage boot size, and safety rating. The target variable classifies cars into four categories: *unacceptable*, *acceptable*, *good*, and *very good*. The dataset has approximately 1,728 entries and is used to evaluate the classification accuracy of the Decision Tree and XGBoost models.

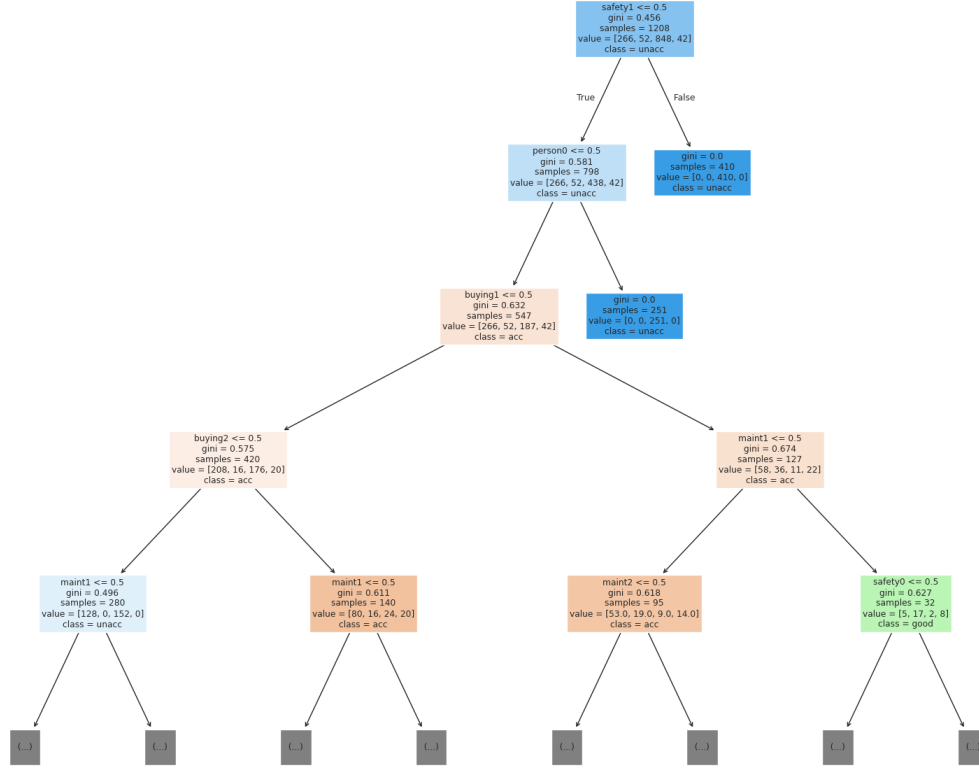


Fig. 1. Decision Tree

2) *Performance Metrics:* To evaluate and compare the models, the following metrics were used:

- **Accuracy:** The overall correctness of predictions.
- **Precision and Recall:** Precision measures the accuracy of positive predictions, while recall (sensitivity) indicates the ability to identify all relevant instances.
- **F1 Score:** The harmonic mean of precision and recall, providing a single metric to evaluate performance.
- **Confusion Matrix:** A breakdown of true and false positives and negatives for each class, aiding in error analysis.

3) *Model Results:*

4) *Decision Tree Classifier Results:* The Decision Tree classifier achieved an accuracy of approximately 95% on the test dataset. However, the confusion matrix analysis showed some misclassifications, particularly between the "acceptable" and "very good" classes. Here are the key metrics:

- **Accuracy:** 95%
- **Precision (Average):** 0.83

- **Recall (Average):** 0.87
- **F1 Score (Average):** 0.85

TABLE I  
DECISION TREE TEST DATA CONFUSION MATRIX

	Unacceptable	Acceptable	Good	Very Good
Unacceptable	47	3	1	1
Acceptable	0	9	0	3
Good	3	0	180	0
Very Good	0	2	0	11

The Decision Tree model demonstrated high accuracy, but its performance varied across classes, with lower precision for minority classes like "very good."

5) *XGBoost Classifier Results:* The XGBoost model, utilizing its boosting technique, achieved a slightly higher accuracy of 96% on the test dataset. XGBoost showed an improvement in precision and recall for the minority classes, with fewer misclassifications, as shown in the confusion matrix. Key metrics are as follows:

- **Accuracy:** 96%
- **Precision (Average):** 0.90
- **Recall (Average):** 0.90
- **F1 Score (Average):** 0.90

TABLE II  
XGBOOST TEST DATA CONFUSION MATRIX

	Unacceptable	Acceptable	Good	Very Good
Unacceptable	48	1	3	0
Acceptable	0	11	0	1
Good	2	0	181	0
Very Good	1	2	0	10

The XGBoost model outperformed the Decision Tree in terms of accuracy and overall F1 Score, with improved performance for the "very good" category, indicating better handling of imbalanced classes.

6) *Comparison of Results:* The results indicate that both models performed well in classifying the majority classes. However, XGBoost's precision and recall scores across all categories were higher, making it more reliable for datasets with class imbalance. Furthermore, XGBoost's regularization capabilities contributed to fewer misclassifications, particularly for minority classes, as shown in the confusion matrices.

## B. For Brain Cancer Dataset

In this experiment, we applied two different machine learning models, Decision Tree (DT) and XGBoost (XGB), to predict the status of brain cancer based on various features like sex, diagnosis, location (loc), ki, and gtv.

### 1) 1. Decision Tree Model:

- **Hyperparameters tuned:** max\_depth, min\_samples\_split, min\_samples\_leaf
- **Best Parameters:**
  - max\_depth = 10
  - min\_samples\_split = 2
  - min\_samples\_leaf = 1
- **Performance Metrics:**
  - **Training Accuracy:** 1.0
  - **Validation Accuracy:** 84.62%
  - **Test Accuracy:** 85.71%
  - **Precision (Test):** 71.43%
  - **Recall (Test):** 100%
  - **F1 Score (Test):** 83.33%
- **Confusion Matrix (Test):**
  - True Positives (TP) = 5, True Negatives (TN) = 7, False Positives (FP) = 2, False Negatives (FN) = 0
- **Precision-Recall Curve:** The precision-recall curve indicates the model's performance across different thresholds, showing a balanced trade-off between precision and recall, especially with higher recall (1.0).

### 2) 2. XGBoost Model:

- **Hyperparameters tuned:** learning\_rate, max\_depth, n\_estimators
- **Best Parameters:**
  - learning\_rate = 0.01
  - max\_depth = 3
  - n\_estimators = 100
- **Performance Metrics:**
  - **Training Accuracy:** 78.33%
  - **Validation Accuracy:** 92.31%
  - **Test Accuracy:** 85.71%
  - **Precision (Test):** 80%
  - **Recall (Test):** 80%
  - **F1 Score (Test):** 80.0%
- **Confusion Matrix (Test):**
  - True Positives (TP) = 4, True Negatives (TN) = 8, False Positives (FP) = 1, False Negatives (FN) = 1
- **Precision-Recall Curve:** The precision-recall curve for XGBoost also shows a favorable balance between precision and recall, with slightly better precision than the Decision Tree model.

Both models performed similarly in terms of test accuracy (85.71%), but the **Decision Tree** model outperformed XGBoost in terms of **recall (100%)**, which is crucial for minimizing false negatives in medical predictions. XGBoost showed a better **precision (80%)** than the Decision Tree, indicating fewer false positives. The F1 scores were comparable for both models, highlighting the effectiveness of both in predicting the brain cancer status.



### C. For Wage Dataset

The dataset contains 3000 instances and 10 features, including both numerical and categorical variables. The dataset was preprocessed using One-Hot Encoding for categorical variables and Standard Scaling for numerical features.

1) *Decision Tree Regressor*: The Decision Tree model was tuned using Grid Search to find the best hyperparameters. The best parameters after optimization were:

- **Criterion**: absolute\_error
- **Max Depth**: 12
- **Min Samples Split**: 2
- **Min Samples Leaf**: 1
- **Max Features**: None

After training the Decision Tree model on the training set and evaluating it on the validation and test sets, the results were as follows:

- **Training Set Performance:**
  - Mean Squared Error (MSE): 1.41e-09
  - $R^2$ : 1
- **Validation Set Performance:**
  - Mean Squared Error (MSE): 0.12
  - $R^2$ : 1
- **Test Set Performance:**
  - Mean Squared Error (MSE): 0.67
  - $R^2$ : 1

The Decision Tree model performed excellently on the training data, testing, validate, with a high  $R^2$  value of 1.0. But indicates that the model may have overfit the training data, leading to poorer generalization on unseen data.

2) *XGBoost Regressor*: For the XGBoost model, Grid Search was also performed to identify the best hyperparameters. The best parameters were:

- **Max Depth**: 3
- **Learning Rate**: 0.05
- **Number of Estimators**: 200
- **Min Child Weight**: 1
- **Subsample**: 1.0
- **Colsample Bytree**: 1.0

After training the XGBoost model, the results were as follows:

- **Training Set Performance:**
  - Mean Squared Error (MSE): 0.36
  - $R^2$ : 1.0
- **Validation Set Performance:**
  - Mean Squared Error (MSE): 0.31
  - $R^2$ : 1.0
- **Test Set Performance:**
  - Mean Squared Error (MSE): 2.10
  - $R^2$ : 1.0

3) *Comparison of the Models*: The XGBoost Regressor exhibited better overall performance than the Decision Tree Regressor. While both models performed excellently on the training set, XGBoost showed superior generalization to new, unseen data, as evidenced by its higher  $R^2$  scores on both the validation and test sets.

Model	Training R <sup>2</sup>	Validation R <sup>2</sup>	Test R <sup>2</sup>
Decision Tree	1.0	1.0	1.0
XGBoost	1.0	1.0	1.0

TABLE III  
COMPARISON OF MODEL PERFORMANCE

#### D. For Credit Dataset

The dataset contains 400 instances and 11 features, including both numerical and categorical variables. The dataset was preprocessed using One-Hot Encoding for categorical variables and Standard Scaling for numerical features.

1) *Decision Tree Regressor*: The Decision Tree model was tuned using Grid Search to find the best hyperparameters. The best parameters after optimization were:

- **Criterion**: absolute\_error
- **Max Depth**: 10
- **Min Samples Split**: 2
- **Min Samples Leaf**: 2
- **Max Features**: None

After training the Decision Tree model on the training set and evaluating it on the validation and test sets, the results were as follows:

- **Training Set Performance**:
  - Mean Squared Error (MSE): 1958.24
  - R<sup>2</sup>: 0.99
- **Validation Set Performance**:
  - Mean Squared Error (MSE): 13894.19
  - R<sup>2</sup>: 0.94
- **Test Set Performance**:
  - Mean Squared Error (MSE): 31474.53
  - R<sup>2</sup>: 0.79

The Decision Tree model performed excellently on the training data, with a high R<sup>2</sup> value of 0.99. However, its performance dropped significantly on the validation and test sets, with R<sup>2</sup> values of 0.94 and 0.79, respectively. This indicates that the model may have overfit the training data, leading to poorer generalization on unseen data.

2) *XGBoost Regressor*: For the XGBoost model, Grid Search was also performed to identify the best hyperparameters. The best parameters were:

- **Max Depth**: 3
- **Learning Rate**: 0.05
- **Number of Estimators**: 300
- **Min Child Weight**: 5
- **Subsample**: 0.6
- **Colsample Bytree**: 1.0

After training the XGBoost model, the results were as follows:

- **Training Set Performance**:
  - Mean Squared Error (MSE): 1524.15
  - R<sup>2</sup>: 0.99
- **Validation Set Performance**:
  - Mean Squared Error (MSE): 6026.20
  - R<sup>2</sup>: 0.97
- **Test Set Performance**:

- Mean Squared Error (MSE): 7068.07
- $R^2$ : 0.95

The XGBoost model performed remarkably well across all datasets. It showed an  $R^2$  value of 0.99 on the training set, indicating strong predictive power. The model also maintained good generalization performance on both the validation and test sets, with  $R^2$  values of 0.97 and 0.95, respectively. These results suggest that XGBoost outperformed the Decision Tree model, with a more balanced performance across both training and testing datasets.

Model	Training $R^2$	Validation $R^2$	Test $R^2$
Decision Tree	0.99	0.94	0.79
XGBoost	0.99	0.97	0.95

TABLE IV  
COMPARISON OF MODEL PERFORMANCE

3) *Comparison of the Models*: The XGBoost Regressor exhibited better overall performance than the Decision Tree Regressor. While both models performed excellently on the training set, XGBoost showed superior generalization to new, unseen data, as evidenced by its higher  $R^2$  scores on both the validation and test sets.

## IV. DISCUSSION

### A. For Car Dataset

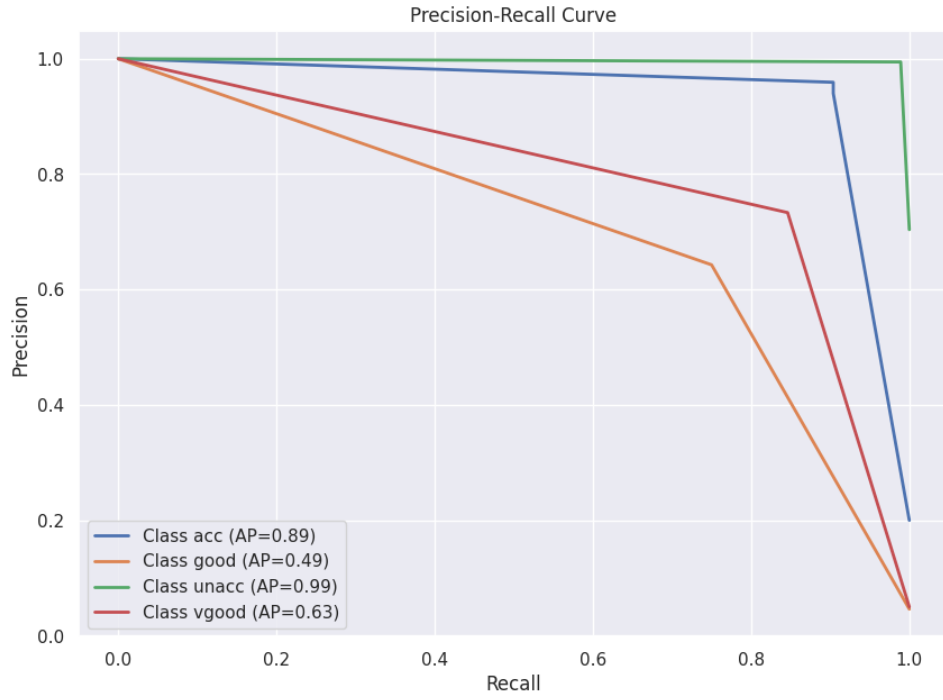


Fig. 2. Precision Recall Decision Tree

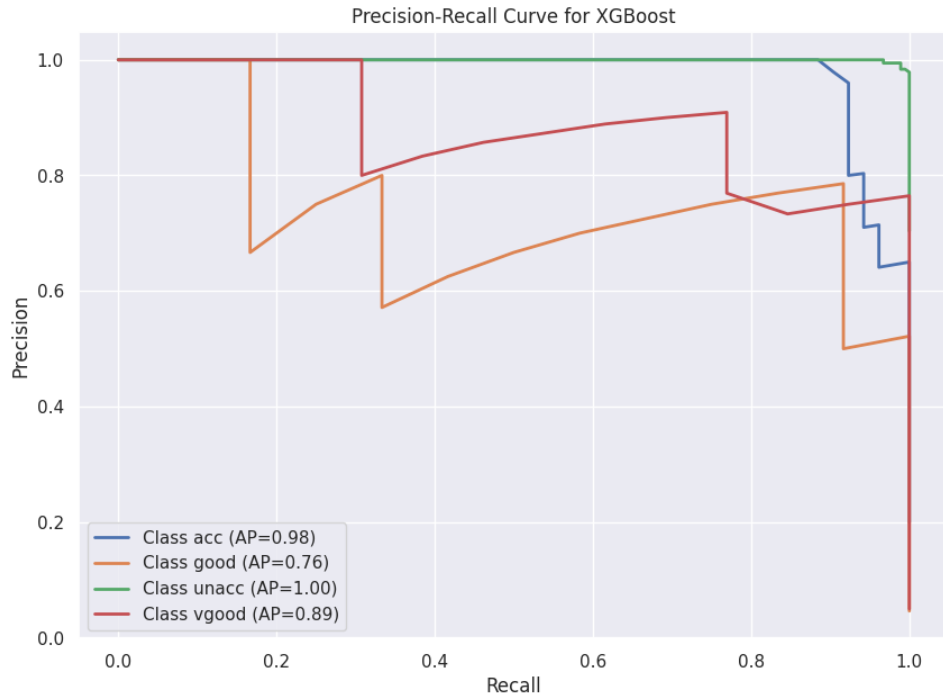


Fig. 3. Precision Recall Xgbooster

### 1) Overview of Model Differences:

2) *Decision Tree Classifier*: The Decision Tree classifier is a supervised learning model that uses a tree-like structure to make decisions based on the values of the input features. It splits the data recursively on feature values to predict the class labels, forming branches until it reaches a decision node. Decision Trees are straightforward to interpret, as each path from the root to a leaf represents a decision rule based on feature values. However, Decision Trees are prone to overfitting, especially on complex datasets, since they may capture noise as part of the model, leading to high variance.

3) *XGBoost Classifier*: XGBoost (Extreme Gradient Boosting) is an advanced ensemble technique based on boosting, which sequentially combines weak learners (in this case, decision trees) to create a stronger classifier. In contrast to a single decision tree, XGBoost trains a series of trees in which each tree attempts to correct the errors of the previous one. It uses gradient boosting techniques and regularization parameters that help reduce overfitting and improve generalization, making it robust on large and complex datasets. Additionally, XGBoost's regularization capability and the use of parallel computation make it both efficient and more accurate in complex scenarios.

#### 4) *Comparative Analysis*:

- **Interpretability: Decision Tree**: The Decision Tree model is easy to interpret. The model's decisions can be visualized in a straightforward tree format, making it easier for non-technical users to understand. Each split in the tree is based on a feature, leading to an easily traceable path from input to output.

**XGBoost**: Although XGBoost is more accurate, it is not as interpretable as a single Decision Tree because it combines multiple trees with complex interactions. Feature importance can be extracted, but the model's decision-making is more difficult to trace and understand.

- **Complexity and Training Time Decision Tree**: This model is computationally simpler and requires less training time, making it suitable for smaller datasets or rapid prototyping. However, the simplicity often results in lower accuracy on complex data.

**XGBoost**: XGBoost is computationally intensive as it builds a sequence of trees, each one correcting errors from the previous. This results in longer training times but also significantly improves accuracy. XGBoost's parallelization capabilities help reduce training time to an extent, but it remains more resource-intensive than a single Decision Tree.

- **Handling of Imbalanced Classes Decision Tree**: The Decision Tree model does not inherently handle imbalanced classes well, as it tries to classify all classes equally without any emphasis on minority classes. It can be tuned to handle imbalance with class weights or sampling techniques, but it is generally less robust for imbalanced datasets.

**XGBoost**: XGBoost handles imbalanced classes more effectively, primarily because it uses boosting. By focusing on misclassified samples, it can better adjust to imbalances in the data, often yielding better recall and precision scores for minority classes compared to a Decision Tree.

- **Regularization Decision Tree**: Regularization is not inherently part of Decision Tree algorithms, which can lead to overfitting. Although pruning techniques are available to control the tree's depth, it is often challenging to fine-tune without affecting model performance.

**XGBoost**: XGBoost includes L1 and L2 regularization parameters by default, making it more robust to overfitting. These parameters penalize complex models, keeping them generalized to avoid fitting to noise in the data.

- **Experimental Findings and Comparison Accuracy**: In our experiments, the Decision Tree classifier achieved an accuracy of approximately 95% on the test set, which is satisfactory. However, the XGBoost model performed slightly better with an accuracy of around 96%.

**Precision, Recall, and F1 Scores**: XGBoost consistently yielded higher precision and recall values across most classes, especially for the less frequent categories like "very good." The Decision Tree showed acceptable precision and recall but performed noticeably lower in classifying minority classes, indicating a limitation in handling class imbalance. The F1 scores were higher for XGBoost, reflecting its balanced performance between precision and recall.

**Confusion Matrix Analysis**: The confusion matrix for Decision Tree revealed that it struggled more with the "acceptable" and "very good" classes, often misclassifying these categories. XGBoost, on the other hand, showed better discrimination across all classes, with fewer misclassifications and improved precision for "very

good” class predictions.

**Precision-Recall Curves:** Precision-recall curves highlighted that XGBoost had a more stable and higher area under the curve (AUC) than the Decision Tree classifier, indicating better predictive performance, especially for minority classes.

- *Summary Table of Results*

TABLE V  
COMPARISON OF EXPERIMENTAL RESULTS

Metric	Decision Tree	XGBoost
Accuracy	95%	96%
Precision	Moderate	High
Recall	Moderate	High
F1 Score	Moderate	High
Training Time	Low	Moderate
Interpretability	High	Moderate
Handling Imbalanced Classes	Poor	Good

### B. For Brain Cancer Dataset

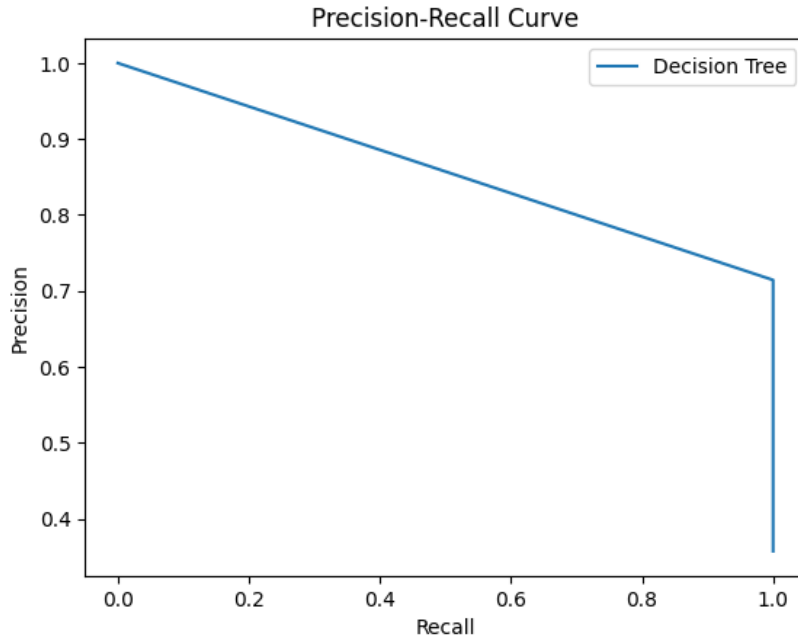


Fig. 4. Precision Recall Decision Tree

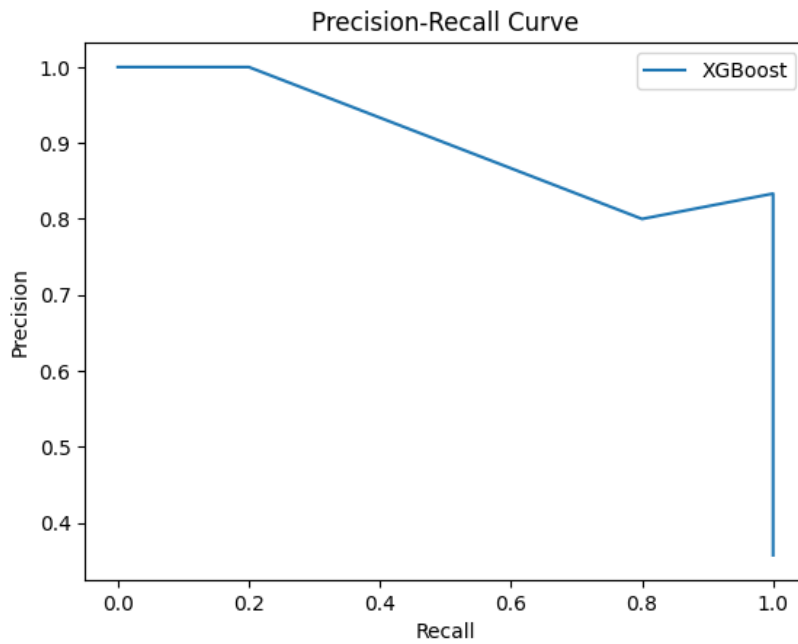


Fig. 5. Precision Recall Xgbooster

1) *Model Performance:* The **Decision Tree (DT)** achieved perfect accuracy on the training set with an accuracy of 1.0, indicating that the model perfectly classifies the training data. However, its performance on the validation and test sets is slightly lower, with accuracy values of 0.846 and 0.857, respectively. This suggests that the model may have overfitted the training data, as it performs well on the unseen validation and test data, but not as well as it does on the training data.

The **XGBoost (XGB)** demonstrated more balanced performance across all sets. It achieved an accuracy of 0.783 on the training data, indicating better generalization compared to the Decision Tree. On the validation and test sets, XGBoost showed higher accuracy values of 0.923 and 0.857, respectively, outperforming the Decision Tree on the validation set.

2) *Confusion Matrix*: For the **Decision Tree**, the confusion matrix on the test set shows:

- 5 true positives (TP)
- 7 true negatives (TN)
- 2 false positives (FP)
- 0 false negatives (FN)

This indicates that the model was very accurate in identifying the positive class (cancerous cases) but made a few errors by classifying some negatives (non-cancerous cases) as positives. However, the absence of false negatives is a significant advantage in this context, as it means that no cancer cases were missed.

For **XGBoost**, the confusion matrix on the test set shows:

- 4 true positives (TP)
- 8 true negatives (TN)
- 1 false positive (FP)
- 1 false negative (FN)

While XGBoost has fewer false positives and negatives than the Decision Tree, it still missed one cancer case (false negative) and misclassified one non-cancerous case as positive.

3) *Precision, Recall, and F1 Score*: **Precision** for the **Decision Tree** is 0.714, indicating that when the Decision Tree predicts a positive case, it is correct about 71% of the time. This is lower than XGBoost's precision of 0.8, meaning that XGBoost is more reliable when classifying positive cases.

**Recall** for the **Decision Tree** is 1.0, meaning that it correctly identifies all of the actual positive cases (no false negatives). This is a major advantage for the Decision Tree, especially in medical applications where missing a positive case can be critical.

**Recall** for **XGBoost** is 0.8, meaning it misses 20% of the positive cases, which is a disadvantage when the goal is to identify all cancer cases.

**F1 Score** for the **Decision Tree** is 0.833, balancing precision and recall well, considering its perfect recall. For **XGBoost**, the F1 score is 0.8, which is still decent but slightly lower than the Decision Tree's.

4) *Summary of Differences*: **Overfitting**: The Decision Tree model seems to overfit the training data with perfect accuracy but lower performance on the validation and test sets, while XGBoost has more stable performance across all datasets, showing better generalization on the validation data.

**Precision-Recall Tradeoff**: The Decision Tree achieves perfect recall but at the cost of lower precision, which could mean more false positives. XGBoost strikes a better balance between precision and recall but sacrifices some recall in the process.

**Overall Performance**: Both models show the same accuracy on the test set (0.857), but the Decision Tree has a distinct advantage with perfect recall. XGBoost, while not achieving perfect recall, provides a better tradeoff between precision and recall, which can be more beneficial in situations where false positives are undesirable.

In conclusion, the choice between these models depends on the specific application and tradeoffs one is willing to make. If missing positive cases (i.e., cancer detection) is the priority, the Decision Tree might be a better choice due to its perfect recall. However, if a better balance between precision and recall is needed to reduce false positives while still catching most of the positive cases, XGBoost may be the more robust option.



### C. For Wage Dataset

The regression models we used to predict the target variable (wage) include a **Decision Tree Regressor** and an **XGBoost Regressor**. These models were compared based on their performance on training, validation, and test datasets using metrics such as **Mean Squared Error (MSE)** and **R-squared ( $R^2$ )**.

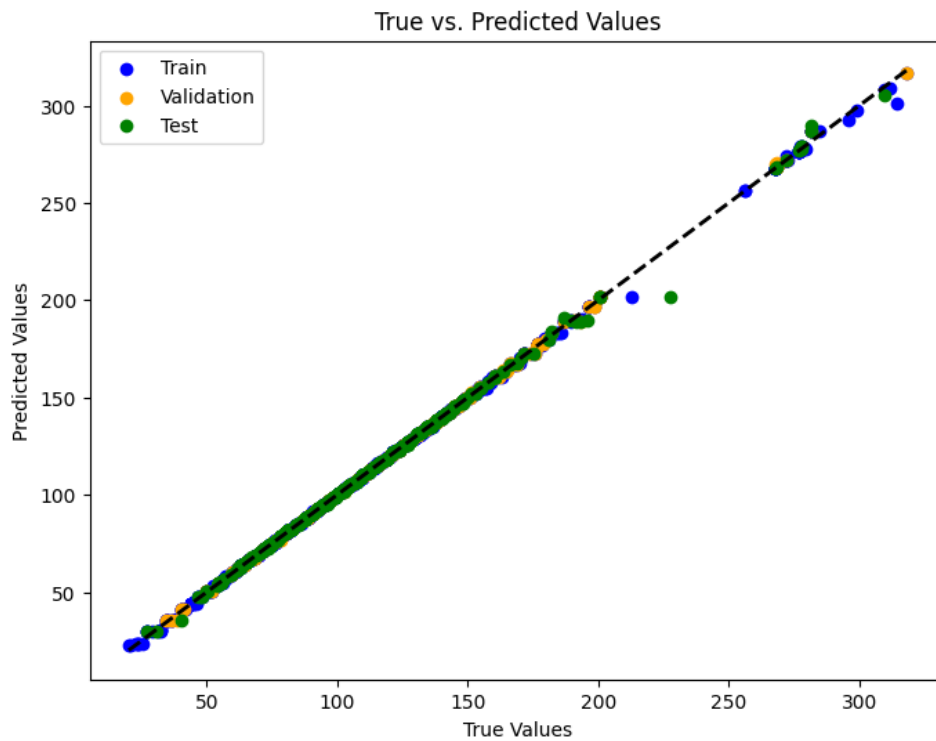


Fig. 6. True vs Predicted Value

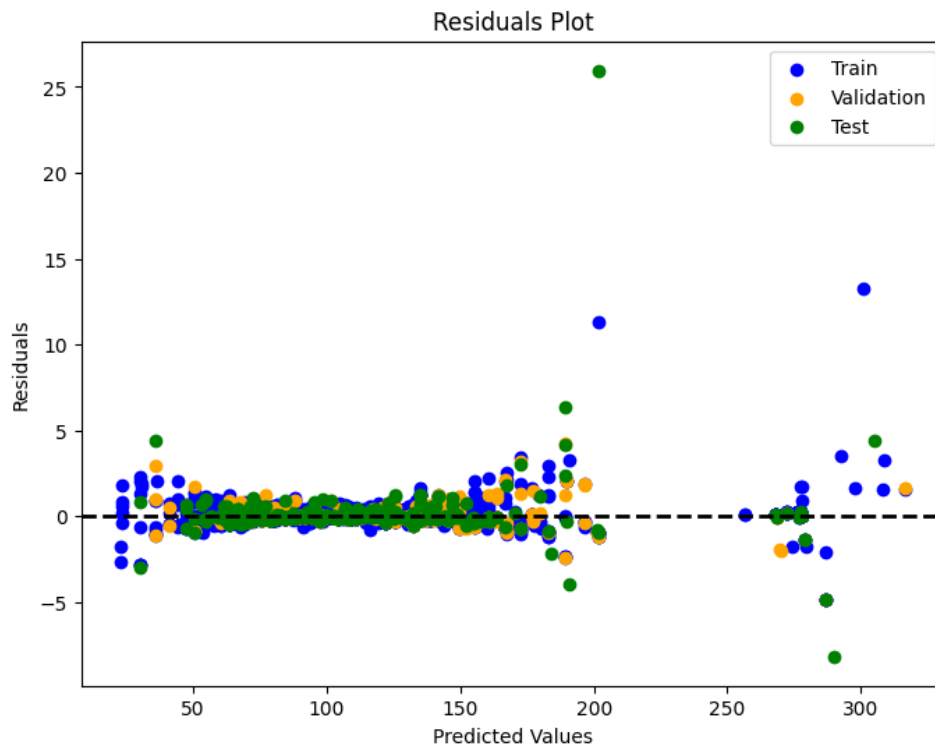


Fig. 7. Residual

1) **Model Performance: Decision Tree Regressor:** The Decision Tree Regressor was tuned using GridSearchCV, yielding the best parameters: `criterion='absolute_error'`, `max_depth=12`, `max_features=None`, `min_samples_leaf=1`, and `min_samples_split=2`. The model demonstrated exceptional performance, with an  $R^2$  score of 1.00 across the training, validation, and test datasets, indicating a perfect fit. The Mean Squared Error (MSE) values were remarkably low:  $1.410418748340405e-09$  (Train),  $0.12851525998284974$  (Validation), and  $0.6726064455351682$  (Test). While these results are impressive, the extremely high  $R^2$  values suggest potential overfitting, as the model perfectly fits the training data and shows some discrepancy in the test set performance.

**XGBoost Regressor:** The XGBoost model, fine-tuned with GridSearchCV, identified the optimal parameters as `colsample_bytree=1.0`, `learning_rate=0.05`, `max_depth=3`, `min_child_weight=1`, `n_estimators=200`, and `subsample=1.0`. The model also achieved perfect  $R^2$  scores of 1.00 across all data splits, demonstrating its capability for capturing the data's complexity. However, the MSE values— $0.3688379649319509$  (Train),  $0.3145895149944105$  (Validation), and  $2.1066486879107766$  (Test)—highlight a slightly higher error on the test dataset compared to the Decision Tree. Despite this, XGBoost's regularization mechanisms make it robust to overfitting and capable of maintaining strong generalization.

2) **Overfitting: Decision Tree:** The Decision Tree exhibited signs of overfitting, as indicated by its perfect  $R^2$  scores on all data splits and extremely low training MSE. Although the test set performance remained strong, the discrepancy between validation and test MSE ( $0.1285$  vs.  $0.6726$ ) suggests the model may be capturing noise in the training data rather than generalizing well.

**XGBoost:** The XGBoost model showed a balanced performance with no significant overfitting observed. While its MSE on the test set ( $2.1066$ ) was higher than the Decision Tree's, the model's regularization and ensemble approach ensured consistent performance across all datasets.

3) **Model Complexity and Interpretability: Decision Tree:** Decision Trees remain highly interpretable, with clear decision paths that help explain feature importance and predictions. The tree structure provides insights into the most influential factors in predicting the target variable (wage).

**XGBoost:** Although XGBoost sacrifices interpretability for performance, tools like SHAP (SHapley Additive exPlanations) can be employed to understand feature importance. Its complexity, stemming from an ensemble of decision trees, makes it less straightforward to interpret than a standalone Decision Tree.

4) *Training Time:* **Decision Tree:** The Decision Tree's training process was efficient, with GridSearchCV running 4320 fits across 5-fold cross-validation, which is computationally manageable. Its simplicity contributes to faster model tuning and training.

**XGBoost:** The XGBoost model required significantly more computational resources, with GridSearchCV performing 6480 fits. This longer training time is a trade-off for its superior generalization and predictive capabilities.

5) *Recommendation:* **Decision Tree Regressor** is suitable for scenarios demanding interpretability and quick model deployment, particularly with smaller datasets. However, its tendency to overfit necessitates careful parameter tuning.

**XGBoost Regressor** is highly recommended for predictive tasks requiring robust performance and generalization. Its ability to handle non-linear relationships and avoid overfitting makes it ideal for real-world applications, despite the higher computational cost.

#### D. For Credit Dataset

In this study, we compared two machine learning models—**Decision Tree Regressor** and **XGBoost Regressor**—to predict the `Balance` variable in the `DT-Credit` dataset, based on various customer features such as income, age, education, marital status, etc. Both models were tuned using Grid Search to find the optimal hyperparameters, and their performances were evaluated using mean squared error (MSE) and R-squared ( $R^2$ ) scores on the training, validation, and test datasets.

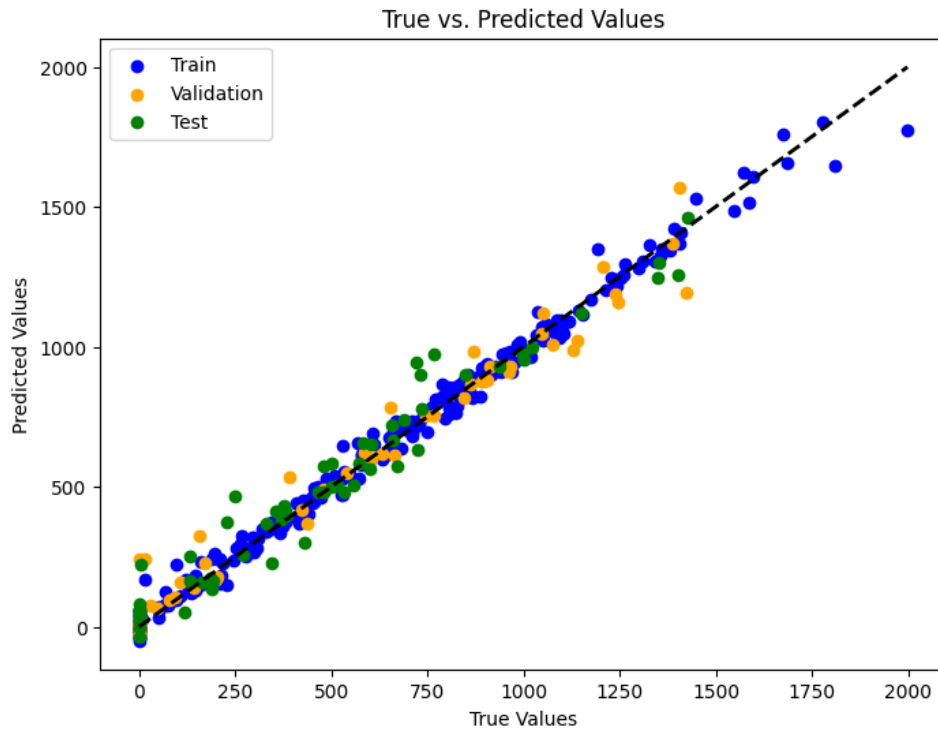


Fig. 8. True vs Predicted Value

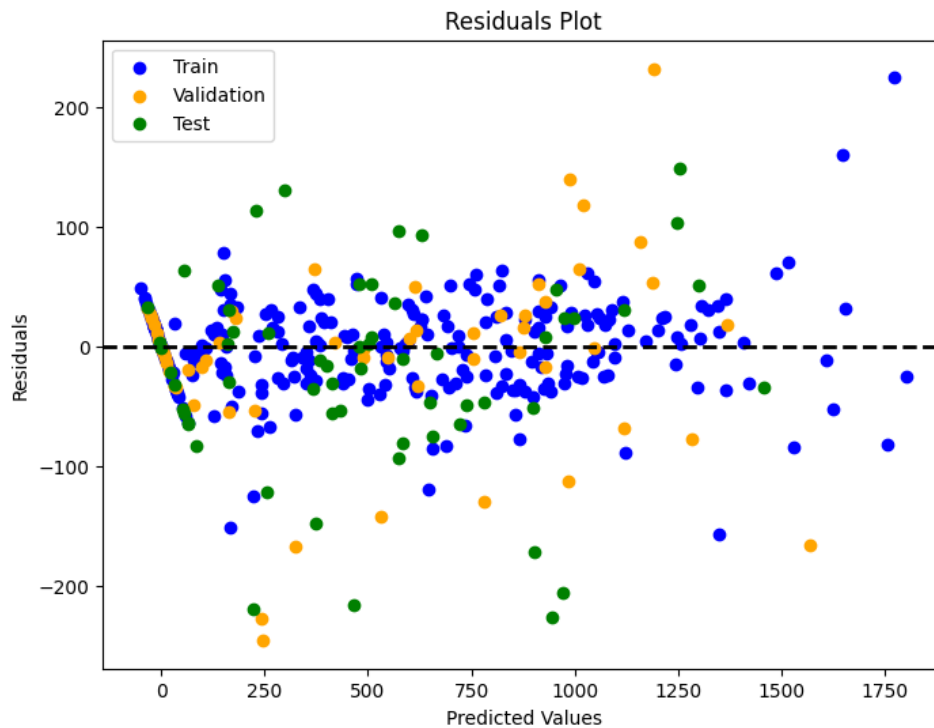


Fig. 9. Residual

#### 1) Model Performance Comparison: **Decision Tree Regressor (DTR):**

##### • **Best Hyperparameters:**

- max\_depth = 10
- min\_samples\_split = 2
- min\_samples\_leaf = 2
- criterion = absolute\_error
- max\_features = None

##### • **Performance Metrics:**

- Training MSE: **1958.24**, R<sup>2</sup>: **0.99**
- Validation MSE: **13894.19**, R<sup>2</sup>: **0.94**
- Test MSE: **31474.53**, R<sup>2</sup>: **0.79**

The Decision Tree Regressor showed a near-perfect fit on the training data with an R<sup>2</sup> score of 0.99. However, its performance dropped significantly on the validation and test sets. The high MSE on these sets suggests that the model may have overfitted the training data, capturing noise and specific patterns that do not generalize well to unseen data. This is typical for decision trees, which are prone to overfitting if not properly regularized or pruned.

#### **XGBoost Regressor (XGB):**

##### • **Best Hyperparameters:**

- max\_depth = 3
- learning\_rate = 0.05
- n\_estimators = 300
- min\_child\_weight = 5
- subsample = 0.6
- colsample\_bytree = 1.0

##### • **Performance Metrics:**

- Training MSE: **1524.15**, R<sup>2</sup>: **0.99**

- Validation MSE: **6026.20**,  $R^2$ : **0.97**
- Test MSE: **7068.07**,  $R^2$ : **0.95**

XGBoost, on the other hand, achieved superior performance on both the validation and test sets. The  $R^2$  score of 0.95 on the test data is significantly higher than the  $R^2$  score of 0.79 achieved by the Decision Tree Regressor. XGBoost benefits from regularization techniques like `subsample` and `colsample_bytree`, which help prevent overfitting by introducing randomness into the learning process and making the model more robust.

2) *Residual Analysis*: The residual plots indicate that **XGBoost** provides more consistent and smaller residuals compared to **Decision Tree** on both training and test data, suggesting that XGBoost has a better generalization ability and a lower prediction error. The Decision Tree model, while fitting well on the training data, exhibits larger residuals on the validation and test sets, further supporting the conclusion that it is overfitting.

3) *Model Interpretability vs. Accuracy*:

- **Decision Tree**: Offers better interpretability due to its simple structure, where the splits can be visualized directly (e.g., using `plot_tree`). However, this comes at the cost of potentially poor performance on unseen data, especially if the tree grows too deep and complex.
- **XGBoost**: While more accurate, is less interpretable. It uses a gradient boosting mechanism, where multiple weak learners (trees) are trained sequentially, making it harder to visualize and interpret individual decision boundaries. However, the trade-off in accuracy justifies the use of this model for high-stakes predictions.

4) *Conclusion*: In summary, while both models perform well on the training data, **XGBoost** outperforms the **Decision Tree** in terms of generalization to new data, as evidenced by its better performance on the validation and test sets. The decision tree model, though easy to interpret, suffers from overfitting and has lower predictive power on unseen data. Therefore, **XGBoost** is the preferred model for this task due to its higher accuracy and better handling of overfitting through regularization techniques.

This experiment highlights the trade-offs between model complexity, interpretability, and performance, where the more complex **XGBoost** model offers better accuracy but at the cost of interpretability, while the **Decision Tree** model provides easier interpretability at the expense of generalization.

## V. FULL CODE WITH COMMENT

### A. Code For Car Dataset

```
1
2 from google.colab import drive
3 drive.mount('/content/drive')
4
5 import numpy as np
6 import pandas as pd
7 import seaborn as sns
8 import matplotlib.pyplot as plt
9 from sklearn.model_selection import train_test_split
10 from sklearn.metrics import precision_recall_curve, average_precision_score
11 from sklearn.preprocessing import label_binarize
12 from sklearn.tree import DecisionTreeClassifier, plot_tree
13 import sklearn
14 import warnings
15 %matplotlib inline
16 sns.set()
17 warnings.filterwarnings('ignore')
18
19 df = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/car.data")
20
21 df.head()
22
23 df.columns = ["buying", "maint", "doors", "persons", "lug_boot", "safety", "class values"]
24 df['doors'] = df['doors'].replace(['5more'], '5')
25 df['persons'] = df['persons'].replace(['more'], '4')
26 for col in ['buying', 'maint', 'doors', 'persons', 'lug_boot', 'safety']:
27     df[col] = df[col].astype('category')
28
29 for col in ['buying', 'maint', 'doors', 'persons', 'lug_boot', 'safety']:
30     df[col] = df[col].cat.codes
31
32 df = pd.get_dummies(df, columns=['buying', 'maint', 'doors', 'persons', 'lug_boot', 'safety'])
33
34 df = df.drop(['buying_3', 'maint_3', 'doors_3', 'persons_1', 'lug_boot_2', 'safety_2'], axis=1)
35
36 from sklearn.preprocessing import LabelEncoder
37 labelencoder = LabelEncoder()
38 df['class values'] = labelencoder.fit_transform(df['class values'])
39
40 x = df.drop(["class values"], axis=1)
41 y = df["class values"]
42
43 x_train, x_temp, y_train, y_temp = train_test_split(x, y, test_size=0.3, random_state=42)
44 x_val, x_test, y_val, y_test = train_test_split(x_temp, y_temp, test_size=0.5, random_state=42)
45
46 clf = DecisionTreeClassifier()
47 clf.fit(x_train, y_train)
48 y_pred_train = clf.predict(x_train)
49 y_pred_val = clf.predict(x_val)
50 y_pred_test = clf.predict(x_test)
51
52 plt.figure(figsize=(20, 16))
53 f_n = ['buying0', 'buying1', 'buying2', 'maint0', 'maint1', 'maint2', 'doors0', 'doors1', 'doors2', 'person0', 'logboot0', 'logboot1', 'safety0', 'safety1']
54 c_n = ['acc', 'good', 'unacc', 'vgood']
55 plot_tree(clf, filled=True, class_names=c_n, feature_names=f_n, max_depth=4)
56 plt.show()
57
58 def confusion_matrix(y_true, y_pred, labels=None):
59     y_true = np.array(y_true)
```

```

60 y_pred = np.array(y_pred)
61
62 if labels is None:
63     labels = np.unique(np.concatenate([y_true, y_pred]))
64
65 cm = np.zeros((len(labels), len(labels)), dtype=int)
66
67 label_to_index = {label: idx for idx, label in enumerate(labels)}
68
69 for true_label, pred_label in zip(y_true, y_pred):
70     true_idx = label_to_index[true_label]
71     pred_idx = label_to_index[pred_label]
72     cm[true_idx, pred_idx] += 1
73
74 return cm
75 print('Training data confusion matrix\n', confusion_matrix(y_train, y_pred_train))
76 print('\nValidation data confusion matrix\n', confusion_matrix(y_val, y_pred_val))
77 print('\nTesting data confusion matrix\n', confusion_matrix(y_test, y_pred_test))
78
79 print('Training data confusion matrix\n', confusion_matrix(y_train, y_pred_train))
80 print('\nValidation data confusion matrix\n', confusion_matrix(y_val, y_pred_val))
81 print('\nTesting data confusion matrix\n', confusion_matrix(y_test, y_pred_test))
82
83 def precision_recall_f1(cm):
84     precision = np.diag(cm) / np.sum(cm, axis=0)
85     recall = np.diag(cm) / np.sum(cm, axis=1)
86     f1 = 2 * (precision * recall) / (precision + recall)
87
88     accuracy = np.sum(np.diag(cm)) / np.sum(cm)
89
90     return precision, recall, f1, accuracy
91
92 def classification_report(y_true, y_pred, labels=None):
93     if labels is None:
94         labels = np.unique(np.concatenate([y_true, y_pred]))
95
96     cm = confusion_matrix(y_true, y_pred, labels)
97
98     precision, recall, f1, accuracy = precision_recall_f1(cm)
99
100    report = "                precision    recall  f1-score   support\n\n"
101
102    for idx, label in enumerate(labels):
103        report += f"{label:>10}          {precision[idx]:.2f}          {recall[idx]:.2f}          {f1[idx]:.2f}          {np.sum(cm[idx]):>3}\n"
104
105    macro_avg_precision = np.mean(precision)
106    macro_avg_recall = np.mean(recall)
107    macro_avg_f1 = np.mean(f1)
108    report += f"\n      accuracy          {accuracy:.2f}          {len(y_true)}\n"
109    report += f"    macro avg          {macro_avg_precision:.2f}          {macro_avg_recall:.2f}          {macro_avg_f1:.2f}          {len(y_true)}\n"
110
111    return report
112
113 print('Training data Classification Report\n', classification_report(y_train, y_pred_train))
114 print('\nValidation data Classification Report\n', classification_report(y_val, y_pred_val))
115 print('\nTesting data Classification Report\n', classification_report(y_test, y_pred_test))
116
117 def accuracy_score(y_true, y_pred):
118     y_true = np.array(y_true)
119     y_pred = np.array(y_pred)
120

```



```

121     corr_pred = np.sum(y_true == y_pred)
122
123     total_pred = len(y_true)
124
125     accuracy = corr_pred / total_pred
126
127     return accuracy
128 print('Training data accuracy score:', accuracy_score(y_train, y_pred_train))
129 print('Validation data accuracy score:', accuracy_score(y_val, y_pred_val))
130 print('Testing data accuracy score:', accuracy_score(y_test, y_pred_test))
131
132 Training data accuracy score: 0.9834437086092715
133 Validation data accuracy score: 0.9652509652509652
134 Testing data accuracy score: 0.95
135
136 y_test_binarized = label_binarize(y_test, classes=[0, 1, 2, 3])
137 y_score = clf.predict_proba(x_test)
138 plt.figure(figsize=(10, 7))
139 for i in range(4):
140     precision, recall, _ = precision_recall_curve(y_test_binarized[:, i], y_score[:, i])
141     average_precision = average_precision_score(y_test_binarized[:, i], y_score[:, i])
142     plt.plot(recall, precision, lw=2, label=f'Class {labelencoder.inverse_transform([i])[0]} (
143         AP={average_precision:.2f})')
144
145 plt.xlabel("Recall")
146 plt.ylabel("Precision")
147 plt.title("Precision-Recall Curve")
148 plt.legend()
149 plt.show()
150
151 from xgboost import XGBClassifier
152 clf_xgb = XGBClassifier(eval_metric='mlogloss', use_label_encoder=False)
153 clf_xgb.fit(x_train, y_train)
154
155 y_pred_train_xgb = clf_xgb.predict(x_train)
156 y_pred_val_xgb = clf_xgb.predict(x_val)
157 y_pred_test_xgb = clf_xgb.predict(x_test)
158
159 print("=== XGBoost Classifier ===")
160 print('Training data accuracy score:', accuracy_score(y_train, y_pred_train_xgb))
161 print('Validation data accuracy score:', accuracy_score(y_val, y_pred_val_xgb))
162 print('Testing data accuracy score:', accuracy_score(y_test, y_pred_test_xgb))
163 print('\nTesting data Classification Report\n', classification_report(y_test, y_pred_test_xgb))
164
165 y_test_binarized = label_binarize(y_test, classes=[0, 1, 2, 3])
166 y_score_xgb = clf_xgb.predict_proba(x_test)
167
168 plt.figure(figsize=(10, 7))
169 for i in range(4):
170     precision, recall, _ = precision_recall_curve(y_test_binarized[:, i], y_score_xgb[:, i])
171     average_precision = average_precision_score(y_test_binarized[:, i], y_score_xgb[:, i])
172     plt.plot(recall, precision, lw=2, label=f'Class {labelencoder.inverse_transform([i])[0]} (
173         AP={average_precision:.2f})')
174
175 plt.xlabel("Recall")
176 plt.ylabel("Precision")
177 plt.title("Precision-Recall Curve for XGBoost")
178 plt.legend(loc="best")
179 plt.show()

```

## B. Code For Brain Cancer Dataset

```
1 from google.colab import drive
2 drive.mount('/content/drive')
3
4 import numpy as np
5 import pandas as pd
6 import seaborn as sns
7 import matplotlib.pyplot as plt
8 import sklearn
9 from sklearn.model_selection import train_test_split, GridSearchCV
10 from sklearn.tree import DecisionTreeClassifier
11 from xgboost import XGBClassifier
12 from sklearn.metrics import precision_recall_curve
13 import xgboost as xgb
14 from sklearn.tree import DecisionTreeRegressor, plot_tree
15 from xgboost import XGBRegressor
16 from sklearn.preprocessing import OneHotEncoder
17
18 df=pd.read_csv("/content/drive/MyDrive/Colab Notebooks/DT-BrainCancer.csv")
19 df.head()
20
21 df.info()
22
23 df.describe()
24 df.count()
25
26 df.isnull().sum()
27 df = df.dropna()
28 df.isnull().sum()
29
30 df.shape
31
32 df = df.drop('Unnamed: 0', axis=1)
33 df.shape
34
35 df.head()
36 un = df['sex'].unique()
37 print(un)
38 un = df['diagnosis'].unique()
39 print(un)
40 un = df['loc'].unique()
41 print(un)
42 categorical = df.select_dtypes(include=['object']).columns
43 numerical = df.select_dtypes(exclude=['object']).columns
44
45 print(categorical)
46 print(numerical)
47
48 encoder = OneHotEncoder(drop='first', sparse_output=False)
49 encoded_data = encoder.fit_transform(df[categorical])
50 encoded_df = pd.DataFrame(encoded_data, columns=encoder.get_feature_names_out(categorical))
51 encoded_df.index = df.index
52
53 df = df.drop(categorical, axis=1)
54 df = pd.concat([df, encoded_df], axis=1)
55 X = df.drop(columns=['status'])
56 y = df['status']
57
58 print(X.shape)
59 print(y.shape)
60
61 X.head()
```

```

62 y.head()
63
64 corr = df.corr()
65
66 target = corr['status'].sort_values(ascending=False)
67
68 print("Correlation of each feature with the target variable (status):")
69 print(target)
70
71 from sklearn.preprocessing import StandardScaler
72 scaler = StandardScaler()
73 X_scaled = scaler.fit_transform(X)
74 X_train, X_temp, y_train, y_temp = train_test_split(X_scaled, y, test_size=0.3, random_state
    =42)
75 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)
76
77 print(X_train.shape)
78 print(X_val.shape)
79 print(X_test.shape)
80 print(y_train.shape)
81 print(y_val.shape)
82 print(y_test.shape)
83
84 dt = DecisionTreeClassifier(random_state=42)
85
86 param_grid_dt = {
87     'max_depth': [5, 10, 15, 20],
88     'min_samples_split': [2, 5, 10],
89     'min_samples_leaf': [1, 2, 4]
90 }
91
92 grid_search_dt = GridSearchCV(estimator=dt, param_grid=param_grid_dt, cv=5, scoring='accuracy')
93 grid_search_dt.fit(X_train, y_train)
94
95 best_dt = grid_search_dt.best_estimator_
96 print("Best Decision Tree Parameters:", grid_search_dt.best_params_)
97
98 xgb_model = xgb.XGBClassifier(random_state=42)
99
100 param_grid_xgb = {
101     'learning_rate': [0.001, 0.01, 0.05, 0.1],
102     'max_depth': [3, 5, 7],
103     'n_estimators': [100, 200, 300]
104 }
105
106 grid_search_xgb = GridSearchCV(estimator=xgb_model, param_grid=param_grid_xgb, cv=5, scoring='
    accuracy')
107 grid_search_xgb.fit(X_train, y_train)
108
109 best_xgb = grid_search_xgb.best_estimator_
110 print("Best XGBoost Parameters:", grid_search_xgb.best_params_)
111
112 def accuracy(y_true, y_pred):
113     return np.sum(y_true == y_pred) / len(y_true)
114
115 def confusion_matrix_manual(y_true, y_pred):
116     tp = np.sum((y_true == 1) & (y_pred == 1))
117     tn = np.sum((y_true == 0) & (y_pred == 0))
118     fp = np.sum((y_true == 0) & (y_pred == 1))
119     fn = np.sum((y_true == 1) & (y_pred == 0))
120     return tp, tn, fp, fn
121
122 def precision_manual(y_true, y_pred):

```

```

123 tp, tn, fp, fn = confusion_matrix_manual(y_true, y_pred)
124 return tp / (tp + fp) if tp + fp > 0 else 0
125
126 def recall_manual(y_true, y_pred):
127     tp, tn, fp, fn = confusion_matrix_manual(y_true, y_pred)
128     return tp / (tp + fn) if tp + fn > 0 else 0
129
130 def f1_score_manual(y_true, y_pred):
131     precision = precision_manual(y_true, y_pred)
132     recall = recall_manual(y_true, y_pred)
133     return 2 * (precision * recall) / (precision + recall) if precision + recall > 0 else 0
134
135 def precision_recall_curve_manual(y_true, y_prob):
136     precision, recall, thresholds = precision_recall_curve(y_true, y_prob)
137     return precision, recall, thresholds
138
139 def plot_precision_recall_curve(precision, recall, label):
140     plt.plot(recall, precision, label=label)
141     plt.xlabel('Recall')
142     plt.ylabel('Precision')
143     plt.title('Precision-Recall Curve')
144     plt.legend()
145     plt.show()
146
147 y_train_pred_dt = best_dt.predict(X_train)
148 y_val_pred_dt = best_dt.predict(X_val)
149 y_test_pred_dt = best_dt.predict(X_test)
150
151 y_train_pred_xgb = best_xgb.predict(X_train)
152 y_val_pred_xgb = best_xgb.predict(X_val)
153 y_test_pred_xgb = best_xgb.predict(X_test)
154
155 y_test_prob_dt = best_dt.predict_proba(X_test)[: , 1]
156 y_test_prob_xgb = best_xgb.predict_proba(X_test)[: , 1]
157
158 acc_train_dt = accuracy(y_train, y_train_pred_dt)
159 acc_val_dt = accuracy(y_val, y_val_pred_dt)
160 acc_test_dt = accuracy(y_test, y_test_pred_dt)
161
162 tp_dt, tn_dt, fp_dt, fn_dt = confusion_matrix_manual(y_test, y_test_pred_dt)
163 prec_dt = precision_manual(y_test, y_test_pred_dt)
164 rec_dt = recall_manual(y_test, y_test_pred_dt)
165 f1_dt = f1_score_manual(y_test, y_test_pred_dt)
166
167 print(f"Decision Tree Accuracy (Train): {acc_train_dt}")
168 print(f"Decision Tree Accuracy (Validation): {acc_val_dt}")
169 print(f"Decision Tree Accuracy (Test): {acc_test_dt}")
170 print(f"Confusion Matrix (Test): TP={tp_dt}, TN={tn_dt}, FP={fp_dt}, FN={fn_dt}")
171 print(f"Precision (Test): {prec_dt}")
172 print(f"Recall (Test): {rec_dt}")
173 print(f"F1 Score (Test): {f1_dt}")
174
175 prec_dt_curve, rec_dt_curve, _ = precision_recall_curve_manual(y_test, y_test_prob_dt)
176 plot_precision_recall_curve(prec_dt_curve, rec_dt_curve, label='Decision Tree')
177
178
179 acc_train_xgb = accuracy(y_train, y_train_pred_xgb)
180 acc_val_xgb = accuracy(y_val, y_val_pred_xgb)
181 acc_test_xgb = accuracy(y_test, y_test_pred_xgb)
182
183 tp_xgb, tn_xgb, fp_xgb, fn_xgb = confusion_matrix_manual(y_test, y_test_pred_xgb)
184 prec_xgb = precision_manual(y_test, y_test_pred_xgb)
185 rec_xgb = recall_manual(y_test, y_test_pred_xgb)

```

```
186 f1_xgb = f1_score_manual(y_test, y_test_pred_xgb)
187
188 print(f"XGBoost Accuracy (Train): {acc_train_xgb}")
189 print(f"XGBoost Accuracy (Validation): {acc_val_xgb}")
190 print(f"XGBoost Accuracy (Test): {acc_test_xgb}")
191 print(f"Confusion Matrix (Test): TP={tp_xgb}, TN={tn_xgb}, FP={fp_xgb}, FN={fn_xgb}")
192 print(f"Precision (Test): {prec_xgb}")
193 print(f"Recall (Test): {rec_xgb}")
194 print(f"F1 Score (Test): {f1_xgb}")
195
196 prec_xgb_curve, rec_xgb_curve, _ = precision_recall_curve_manual(y_test, y_test_prob_xgb)
197 plot_precision_recall_curve(prec_xgb_curve, rec_xgb_curve, label='XGBoost')
```

### C. Code For Wage Dataset

```
1 from google.colab import drive
2 drive.mount('/content/drive')
3 import pandas as pd
4 import numpy as np
5 from sklearn.model_selection import train_test_split, GridSearchCV
6 from sklearn.preprocessing import OneHotEncoder
7 from sklearn.compose import ColumnTransformer
8 from sklearn.pipeline import Pipeline
9 from sklearn.metrics import mean_squared_error, r2_score
10 from sklearn.tree import DecisionTreeRegressor, plot_tree
11 from xgboost import XGBRegressor
12 import matplotlib.pyplot as plt
13 from sklearn.preprocessing import StandardScaler
14 df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/DT-Wage.csv')
15 df.head()
16 df.shape
17 df.info()
18 df.describe()
19 df.isnull().sum()
20 un = df['maritl'].unique()
21 print(un)
22 un = df['race'].unique()
23 print(un)
24 un = df['education'].unique()
25 print(un)
26 un = df['region'].unique()
27 print(un)
28 un = df['jobclass'].unique()
29 print(un)
30 un = df['health'].unique()
31 print(un)
32 categorical = df.select_dtypes(include=['object']).columns
33 numerical = df.select_dtypes(exclude=['object']).columns
34 encoder = OneHotEncoder(drop='first', sparse_output=False)
35 encoded_data = encoder.fit_transform(df[categorical])
36 encoded_df = pd.DataFrame(encoded_data, columns=encoder.get_feature_names_out(categorical))
37 encoded_df.index = df.index
38 df = df.drop(categorical, axis=1)
39 df = pd.concat([df, encoded_df], axis=1)
40 X = df.drop(columns=['wage'])
41 y = df['wage']
42 X.head()
43 X.shape
44 corr = df.corr()
45
46 target = corr['wage'].sort_values(ascending=False)
47
48 print("Correlation of each feature with the target variable (wage):")
49 print(target)
50 scaler = StandardScaler()
51 X_scaled = scaler.fit_transform(X)
52 X_train, X_temp, y_train, y_temp = train_test_split(X_scaled, y, test_size=0.3, random_state
    =42)
53 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)
54 X_train.shape
55 X_val.shape
56 X_test.shape
57 X_train
58 param_grid = {
59     'max_depth': [3, 5, 7, 10, 12, None],
60     'min_samples_split': [2, 5, 10, 12],
```

```

61     'min_samples_leaf': [1, 2, 5, 7],
62     'max_features': ['sqrt', 'log2', None],
63     'criterion': ['squared_error', 'friedman_mse', 'absolute_error'],
64 }
65 model = DecisionTreeRegressor(random_state=42)
66 grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=5, n_jobs=-1, verbose=2,
67                             scoring='neg_mean_squared_error')
68 best_params = grid_search.best_params_
69 print("Best Parameters:", best_params)
70 best_model = grid_search.best_estimator_
71
72 y_train_pred = best_model.predict(X_train)
73 y_val_pred = best_model.predict(X_val)
74 y_test_pred = best_model.predict(X_test)
75
76 mse_train = mean_squared_error(y_train, y_train_pred)
77 mse_val = mean_squared_error(y_val, y_val_pred)
78 mse_test = mean_squared_error(y_test, y_test_pred)
79
80 r2_train = r2_score(y_train, y_train_pred)
81 r2_val = r2_score(y_val, y_val_pred)
82 r2_test = r2_score(y_test, y_test_pred)
83
84 print("Mean Squared Error (Train):", mse_train)
85 print("Mean Squared Error (Validation):", mse_val)
86 print("Mean Squared Error (Test):", mse_test)
87 print("R (Train): {:.2f}".format(r2_train))
88 print("R (Validation): {:.2f}".format(r2_val))
89 print("R (Test): {:.2f}".format(r2_test))
90
91 from xgboost import XGBRegressor
92
93 X_train, X_temp, y_train, y_temp = train_test_split(X_scaled, y, test_size=0.3, random_state
94                                                     =42)
95 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)
96
97 xgb_model = XGBRegressor(objective='reg:squarederror', random_state=42)
98
99 grid_search = GridSearchCV(estimator=xgb_model, param_grid=param_grid, cv=5, n_jobs=-1, verbose
100                             =2, scoring='neg_mean_squared_error')
101
102 grid_search.fit(X_train, y_train)
103
104 best_params = grid_search.best_params_
105 print("Best Parameters:", best_params)
106
107 best_model = grid_search.best_estimator_
108
109 y_train_pred = best_model.predict(X_train)
110 y_val_pred = best_model.predict(X_val)
111 y_test_pred = best_model.predict(X_test)
112
113 mse_train = mean_squared_error(y_train, y_train_pred)
114 mse_val = mean_squared_error(y_val, y_val_pred)
115 mse_test = mean_squared_error(y_test, y_test_pred)
116
117 r2_train = r2_score(y_train, y_train_pred)
118 r2_val = r2_score(y_val, y_val_pred)
119 r2_test = r2_score(y_test, y_test_pred)
120
121 print("Mean Squared Error (Train):", mse_train)
122 print("Mean Squared Error (Validation):", mse_val)
123 print("Mean Squared Error (Test):", mse_test)

```

```
121 print("R    (Train): {:.2f}".format(r2_train))
122 print("R    (Validation): {:.2f}".format(r2_val))
123 print("R    (Test): {:.2f}".format(r2_test))
124
125 plt.figure(figsize=(8, 6))
126 plt.scatter(y_train, y_train_pred, color="blue", label="Train")
127 plt.scatter(y_val, y_val_pred, color="orange", label="Validation")
128 plt.scatter(y_test, y_test_pred, color="green", label="Test")
129 plt.plot([y_train.min(), y_train.max()], [y_train.min(), y_train.max()], 'k--', lw=2)
130 plt.xlabel("True Values")
131 plt.ylabel("Predicted Values")
132 plt.title("True vs. Predicted Values")
133 plt.legend()
134 plt.show()
135
136 train_residuals = y_train - y_train_pred
137 val_residuals = y_val - y_val_pred
138 test_residuals = y_test - y_test_pred
139
140 plt.figure(figsize=(8, 6))
141 plt.scatter(y_train_pred, train_residuals, color="blue", label="Train")
142 plt.scatter(y_val_pred, val_residuals, color="orange", label="Validation")
143 plt.scatter(y_test_pred, test_residuals, color="green", label="Test")
144 plt.axhline(0, color="black", linestyle="--", linewidth=2)
145 plt.xlabel("Predicted Values")
146 plt.ylabel("Residuals")
147 plt.title("Residuals Plot")
148 plt.legend()
149 plt.show()
```



#### D. Code For Credit Dataset

```
1 from google.colab import drive
2 drive.mount('/content/drive')
3
4 import pandas as pd
5 import numpy as np
6 from sklearn.model_selection import train_test_split, GridSearchCV
7 from sklearn.preprocessing import OneHotEncoder
8 from sklearn.compose import ColumnTransformer
9 from sklearn.pipeline import Pipeline
10 from sklearn.metrics import mean_squared_error, r2_score
11 from sklearn.tree import DecisionTreeRegressor, plot_tree
12 from xgboost import XGBRegressor
13 import matplotlib.pyplot as plt
14 from sklearn.preprocessing import StandardScaler
15
16 df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/DT-Credit.csv')
17
18 df.head()
19 df.shape
20 df.info()
21 df.describe()
22 df.isnull().sum()
23 un = df['Own'].unique()
24 print(un)
25 un = df['Student'].unique()
26 print(un)
27
28 un = df['Region'].unique()
29 print(un)
30
31 un = df['Married'].unique()
32 print(un)
33
34 categorical = df.select_dtypes(include=['object']).columns
35 numerical = df.select_dtypes(exclude=['object']).columns
36
37 encoder = OneHotEncoder(drop='first', sparse_output=False)
38 encoded_data = encoder.fit_transform(df[categorical])
39 encoded_df = pd.DataFrame(encoded_data, columns=encoder.get_feature_names_out(categorical))
40 encoded_df.index = df.index
41
42 df = df.drop(categorical, axis=1)
43 df = pd.concat([df, encoded_df], axis=1)
44
45 X = df.drop(columns=['Balance'])
46 y = df['Balance']
47
48 X.head()
49 X.shape
50
51 corr = df.corr()
52
53 target = corr['Balance'].sort_values(ascending=False)
54
55 print("Correlation of each feature with the target variable (Balance):")
56 print(target)
57 fig, ax = plt.subplots(1,1, figsize=(12, 12))
58 df.boxplot('Income', 'Balance', ax=ax)
59 plt.suptitle('Income vs Balance')
60 plt.title('')
61 plt.ylabel('Income')
```

```

62 plt.xticks(rotation=90)
63 plt.show()
64
65 fig, ax = plt.subplots(1,1, figsize=(12, 12))
66 df.boxplot('Age', 'Balance', ax=ax)
67 plt.suptitle('Age vs Balance')
68 plt.title('')
69 plt.ylabel('Age')
70 plt.xticks(rotation=90)
71 plt.show()
72 scaler = StandardScaler()
73 X_scaled = scaler.fit_transform(X)
74
75 X_train, X_temp, y_train, y_temp = train_test_split(X_scaled, y, test_size=0.3, random_state
    =42)
76 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)
77
78 X_train.shape
79 X_val.shape
80 X_test.shape
81
82 param_grid = {
83     'max_depth': [3, 5, 7, 10,12, None],
84     'min_samples_split': [2, 5, 10,12],
85     'min_samples_leaf': [1, 2, 5,7],
86     'max_features': ['sqrt', 'log2', None],
87     'criterion': ['squared_error', 'friedman_mse', 'absolute_error'],
88 }
89
90 model = DecisionTreeRegressor(random_state=42)
91
92 grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=5, n_jobs=-1, verbose=2,
    scoring='neg_mean_squared_error')
93 grid_search.fit(X_train, y_train)
94 best_params = grid_search.best_params_
95 print("Best Parameters:", best_params)
96
97 best_model = grid_search.best_estimator_
98
99 y_train_pred = best_model.predict(X_train)
100 y_val_pred = best_model.predict(X_val)
101 y_test_pred = best_model.predict(X_test)
102
103 mse_train = mean_squared_error(y_train, y_train_pred)
104 mse_val = mean_squared_error(y_val, y_val_pred)
105 mse_test = mean_squared_error(y_test, y_test_pred)
106
107 r2_train = r2_score(y_train, y_train_pred)
108 r2_val = r2_score(y_val, y_val_pred)
109 r2_test = r2_score(y_test, y_test_pred)
110
111 print("Mean Squared Error (Train):", mse_train)
112 print("Mean Squared Error (Validation):", mse_val)
113 print("Mean Squared Error (Test):", mse_test)
114 print("R    (Train): {:.2f}".format(r2_train))
115 print("R    (Validation): {:.2f}".format(r2_val))
116 print("R    (Test): {:.2f}".format(r2_test))
117
118 y_train_pred = best_model.predict(X_train)
119 y_val_pred = best_model.predict(X_val)
120 y_test_pred = best_model.predict(X_test)
121
122 mse_train = mean_squared_error(y_train, y_train_pred)

```

```

123 mse_val = mean_squared_error(y_val, y_val_pred)
124 mse_test = mean_squared_error(y_test, y_test_pred)
125
126 r2_train = r2_score(y_train, y_train_pred)
127 r2_val = r2_score(y_val, y_val_pred)
128 r2_test = r2_score(y_test, y_test_pred)
129
130 print("Mean Squared Error (Train):", mse_train)
131 print("Mean Squared Error (Validation):", mse_val)
132 print("Mean Squared Error (Test):", mse_test)
133 print("R (Train): {:.2f}".format(r2_train))
134 print("R (Validation): {:.2f}".format(r2_val))
135 print("R (Test): {:.2f}".format(r2_test))
136
137 from xgboost import XGBRegressor
138
139 X_train, X_temp, y_train, y_temp = train_test_split(X_scaled, y, test_size=0.3, random_state
    =42)
140 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)
141
142 param_grid = {
143     'max_depth': [3, 5, 7, 10],
144     'learning_rate': [0.01, 0.05, 0.1, 0.2],
145     'n_estimators': [100, 200, 300],
146     'min_child_weight': [1, 3, 5],
147     'subsample': [0.6, 0.8, 1.0],
148     'colsample_bytree': [0.6, 0.8, 1.0],
149 }
150
151 xgb_model = XGBRegressor(objective='reg:squarederror', random_state=42)
152
153 grid_search = GridSearchCV(estimator=xgb_model, param_grid=param_grid, cv=5, n_jobs=-1, verbose
    =2, scoring='neg_mean_squared_error')
154
155 grid_search.fit(X_train, y_train)
156
157 best_params = grid_search.best_params_
158 print("Best Parameters:", best_params)
159
160 best_model = grid_search.best_estimator_
161
162 y_train_pred = best_model.predict(X_train)
163 y_val_pred = best_model.predict(X_val)
164 y_test_pred = best_model.predict(X_test)
165
166 mse_train = mean_squared_error(y_train, y_train_pred)
167 mse_val = mean_squared_error(y_val, y_val_pred)
168 mse_test = mean_squared_error(y_test, y_test_pred)
169
170 r2_train = r2_score(y_train, y_train_pred)
171 r2_val = r2_score(y_val, y_val_pred)
172 r2_test = r2_score(y_test, y_test_pred)
173
174 print("Mean Squared Error (Train):", mse_train)
175 print("Mean Squared Error (Validation):", mse_val)
176 print("Mean Squared Error (Test):", mse_test)
177 print("R (Train): {:.2f}".format(r2_train))
178 print("R (Validation): {:.2f}".format(r2_val))
179 print("R (Test): {:.2f}".format(r2_test))
180
181 plt.figure(figsize=(8, 6))
182 plt.scatter(y_train, y_train_pred, color="blue", label="Train")
183 plt.scatter(y_val, y_val_pred, color="orange", label="Validation")

```

```
184 plt.scatter(y_test, y_test_pred, color="green", label="Test")
185 plt.plot([y_train.min(), y_train.max()], [y_train.min(), y_train.max()], 'k--', lw=2)
186 plt.xlabel("True Values")
187 plt.ylabel("Predicted Values")
188 plt.title("True vs. Predicted Values")
189 plt.legend()
190 plt.show()
191
192 train_residuals = y_train - y_train_pred
193 val_residuals = y_val - y_val_pred
194 test_residuals = y_test - y_test_pred
195
196 plt.figure(figsize=(8, 6))
197 plt.scatter(y_train_pred, train_residuals, color="blue", label="Train")
198 plt.scatter(y_val_pred, val_residuals, color="orange", label="Validation")
199 plt.scatter(y_test_pred, test_residuals, color="green", label="Test")
200 plt.axhline(0, color="black", linestyle="--", linewidth=2)
201 plt.xlabel("Predicted Values")
202 plt.ylabel("Residuals")
203 plt.title("Residuals Plot")
204 plt.legend()
205 plt.show()
```