# Assignment-5: ZeroR classifier, OneR classifier,K-Nearest-Neighbor Classifiers,Naive Bayesian Classifier, SVM,SVR

Mahbub Ahmed Turza

*ID: 2211063042*

*North South University*

mahbub.turza@northsouth.edu

November 23, 2024

### Abstract

In machine learning, supervised learning is very important as it provides a variety of methods for tasks including regression and classification. In this work, we look into the efficiency of multiple datasets with different properties using traditional methods for supervised learning, such as ZeroR, OneR, K-Nearest Neighbors (K-NN), Naive Bayes, Support Vector Machines (SVM), and Support Vector Regression (SVR). Accuracy, precision, recall, F1-score, precision-recall curves, and confusion matrix are utilized to evaluate performance for binary and multi-class classification problems. Mean squared error (MSE) is a method to evaluate regression performance. We evaluate these algorithms on a range of tasks using thorough experimental analysis, emphasizing their robustness in managing complicated data, computational efficiency, and forecast accuracy. The results show the relative benefits and drawbacks of different approaches and offer helpful suggestions for using them in actual situations.

## I. INTRODUCTION

Supervised learning is vital for predictive analytics, because machine learning is emerging as a revolutionary technology for modern computational research. Using labeled datasets, a model has been trained to infer a mapping from inputs to outputs in the conventional supervised learning paradigm. Because of their interpretability, computational economy, and effectiveness in low-data regimes, classical supervised algorithms remain to be essential to machine learning even in the face of the growth of advanced techniques like deep learning.

Six classic supervised learning techniques are examined in this paper: Support Vector Machines (SVM), Support Vector Regression (SVR), ZeroR, OneR, K-Nearest Neighbors (K-NN), and Naive Bayesian Classifier. To evaluate these algorithms' practicality and generalizability, they were used on a variety of datasets which included regression, multi-class classification, and binary classification problems.

The evaluation methodology was developed to provide an in-depth assessment of algorithmic performance through the use of established indicators. We use precision, recall, F1-scores, accuracy, precision-recall curves, and confusion matrices for binary classification. Mean squared error (MSE) can be utilized to evaluate regression performance, while average accuracy, recall, and F1-score metrics are used to evaluate multi-class classification.

The results of this study provide helpful data for choosing suitable methods in various problem domains by benchmarking these classical algorithms across a variety of tasks and datasets. Additionally, the comparison study draws attention to trade-offs between accuracy, computing cost, and complexity, which enhances the understanding of supervised learning paradigms in both academic and business contexts.

## II. METHODOLOGY

### A. Brain Dataset

*1) Data Collection:* For this study, we used the `BrainCancer` The dataset consists of several features and a target variable representing class labels. The dataset can be loaded and inspected using Pandas.

*2) Data Preprocessing:* Preprocessing is essential to ensure that the data is in an appropriate format for machine learning algorithms. The following steps were implemented for data preprocessing:

*3) Handling Missing Values:* Missing data is handled by filling in or removing rows or columns with missing values.

```python
import pandas as pd
df=pd.read_csv("/content/drive/MyDrive/Colab Notebooks/DT-BrainCancer.csv")
df.fillna(df.mean(), inplace=True)
```

*4) One-Hot Encoding:* One-Hot Encoding is applied to convert categorical variables into numerical form. Each category is transformed into a new binary column, indicating the presence or absence of that category.

```python
df_encoded = pd.get_dummies(df, drop_first=True)
```

One-Hot Encoding is crucial because many algorithms, such as K-NN and SVM, do not support categorical variables directly. These algorithms require the data to be in numerical form, and One-Hot Encoding efficiently transforms categorical data into binary values.

*5) Feature Scaling:* Feature scaling is applied to bring all features into the same scale, especially for algorithms like K-NN and SVM that are sensitive to the scale of the features.

```python
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
df_encoded[['feature1', 'feature2', 'feature3']] = scaler.fit_transform(df_encoded[['feature1',
    'feature2', 'feature3']])
```

*6) Split the Data:* After pre-processing, the dataset is split into training and testing subsets, with the training data used to train the classifiers and the test data used to evaluate model performance.

```python
from sklearn.model_selection import train_test_split

X = df_encoded.drop('target', axis=1)
y = df_encoded['target']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

*7) Classifiers Implementation:* Now, we implement the classifiers for comparison: K-Nearest Neighbors (K-NN), Naive Bayes, Support Vector Machine (SVM), OneR, and ZeroR.

*8) K-Nearest Neighbors (K-NN):* K-NN is a non-parametric method used for classification. The algorithm works by finding the k nearest neighbors to a test instance and assigning the majority class.

```python
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score

knn = KNeighborsClassifier(n_neighbors=5)

knn.fit(X_train, y_train)

y_pred_knn = knn.predict(X_test)

accuracy_knn = accuracy_score(y_test, y_pred_knn)
print(f"K-NN Accuracy: {accuracy_knn}")
```

*9) Naive Bayes:* Naive Bayes classifiers are based on applying Bayes' theorem with strong (naive) independence assumptions. It is effective for text classification and datasets with categorical features.

```python
from sklearn.naive_bayes import GaussianNB

nb = GaussianNB()

nb.fit(X_train, y_train)

y_pred_nb = nb.predict(X_test)

accuracy_nb = accuracy_score(y_test, y_pred_nb)
print(f"Naive Bayes Accuracy: {accuracy_nb}")
```

*10) Support Vector Machine (SVM):* SVM is a supervised learning algorithm that classifies data by finding the hyperplane that best separates the classes.

```python
pipeline = Pipeline([
    ('scaler', StandardScaler()),
    ('svc', SVC(probability=True))
])

# Define the parameter grid
param_grid = {
    'svc__C': [0.01, 0.1, 1, 10, 100],
    'svc__kernel': ['linear', 'rbf'],
    'svc__gamma': ['scale', 0.001, 0.01, 0.1, 1]
}

grid_search = GridSearchCV(pipeline, param_grid, cv=5, scoring='roc_auc', verbose=2)
grid_search.fit(X_train, y_train)

best_model = grid_search.best_estimator_

print("Best Parameters:", grid_search.best_params_)
print("Best ROC-AUC Score:", grid_search.best_score_)
```

*11) OneR Classifier (Implemented Without Library):* OneR (One Rule) is a simple classifier that generates a single rule based on one attribute that provides the best classification performance. It's easy to implement and offers an interpretable model.

```python
import numpy as np

def oneR(X_train, y_train):
    best_accuracy = 0
    best_column = -1
    best_rule = None

    for col in X_train.columns:
        value_counts = X_train[col].value_counts()
        best_value = value_counts.idxmax()
        predictions = (X_train[col] == best_value).astype(int)
        accuracy = np.mean(predictions == y_train)
        if accuracy > best_accuracy:
            best_accuracy = accuracy
            best_column = col
            best_rule = best_value

    return best_column, best_rule

best_column, best_rule = oneR(X_train, y_train)
print(f"Best attribute for OneR: {best_column}, Best rule: {best_rule}")
```

*12) ZeroR Classifier (Implemented Without Library):* ZeroR is the most basic classifier, which simply predicts the majority class for all instances. It is used as a baseline comparison.

```python
def zeroR(X_train, y_train):
    majority_class = y_train.value_counts().idxmax()
    predictions = [majority_class] * len(y_train)

    accuracy = np.mean(predictions == y_train)
    return accuracy

accuracy_zero_r = zeroR(X_train, y_train)
print(f"ZeroR Accuracy: {accuracy_zero_r}")
```

*13) Model Evaluation:* After training the models, we evaluate their performance using the following metrics:

- **Accuracy**: The proportion of correctly classified instances.
- **Precision**: The proportion of true positive predictions among all positive predictions.
- **Recall**: The proportion of true positive predictions among all actual positive instances.
- **F1-Score**: The harmonic mean of precision and recall.
- **Confusion Matrix**: A summary of prediction results showing the true positives, false positives, true negatives, and false negatives.

```python
from sklearn.metrics import classification_report, confusion_matrix

print("Classification Report for K-NN:")
print(classification_report(y_test, y_pred_knn))

print("Confusion Matrix for K-NN:")
print(confusion_matrix(y_test, y_pred_knn))
# Repeat the above for all classifiers
```

*14) Hyperparameter Tuning for K-NN:* For K-NN, hyperparameter tuning is essential to find the optimal number of neighbors (k). We use GridSearchCV to tune the parameter.

```python
from sklearn.model_selection import GridSearchCV

param_grid = {'n_neighbors': [1, 3, 5, 7, 9]}
grid_search_knn = GridSearchCV(KNeighborsClassifier(), param_grid, cv=5)

grid_search_knn.fit(X_train, y_train)

print("Best parameters for K-NN:", grid_search_knn.best_params_)
print("Best score for K-NN:", grid_search_knn.best_score_)
```

## B. Car Dataset

*1) Data Preprocessing:* The dataset used in this study was preprocessed to ensure it was in a format suitable for machine learning models. The preprocessing steps involved handling missing data, scaling features, and encoding categorical variables. The features were standardized to ensure that all variables had the same scale. Categorical features were encoded using Label Encoding to convert them into numerical values.

*2) Model Selection:* We selected the following models for comparison:

- **ZeroR**: A baseline model that predicts the majority class without considering the features.
- **OneR**: A simple model that selects the best single feature to minimize classification error.
- **K-Nearest Neighbors (KNN)**: A non-parametric algorithm that classifies based on the majority vote of the nearest neighbors.
- **Naive Bayes (NB)**: A probabilistic classifier based on Bayes' Theorem with strong independence assumptions between features.
- **Support Vector Machine (SVM)**: A powerful classifier that finds the hyperplane that best separates the data into classes.

All models were trained using the same training data and evaluated based on their performance on the validation and test datasets.

*3) ZeroR (Baseline Model):* ZeroR is the simplest form of classification where the model predicts the most frequent class without considering the features. This model serves as a baseline for comparison with other more complex models. The implementation for ZeroR is as follows:

```
1  from sklearn.dummy import DummyClassifier
2
3  zeroR = DummyClassifier(strategy='most_frequent')
4  zeroR.fit(x_train, y_train)
5  y_val_pred_zeroR = zeroR.predict(x_val)
6  y_test_pred_zeroR = zeroR.predict(x_test)
```

Listing 1: ZeroR Model Implementation

*4) OneR:* OneR is a simple rule-based classifier that selects a single feature and constructs a rule based on this feature to classify the data. The feature with the lowest classification error is selected. OneR can be implemented by evaluating each feature's performance and selecting the best one.

```
1  majority_class = y_train.mode()[0]
2  print(f"Majority class (ZeroR prediction): {majority_class}")
```

Listing 2: OneR Model Implementation

*5) K-Nearest Neighbors (KNN):* KNN is a simple non-parametric method that classifies a data point based on the majority class of its 'k' nearest neighbors. The implementation for KNN is as follows:

```
1  from sklearn.neighbors import KNeighborsClassifier
2
3  knn = KNeighborsClassifier(n_neighbors=5)
4  knn.fit(x_train, y_train)
5  y_val_pred_knn = knn.predict(x_val)
6  y_test_pred_knn = knn.predict(x_test)
```

Listing 3: KNN Model Implementation

*6) Naive Bayes (NB):* Naive Bayes is a probabilistic classifier that applies Bayes' Theorem with the assumption that the features are conditionally independent. The implementation for Naive Bayes is as follows:

```
1  from sklearn.naive_bayes import GaussianNB
2
3  nb = GaussianNB()
4  nb.fit(x_train, y_train)
5  y_val_pred_nb = nb.predict(x_val)
```

```
6  y_test_pred_nb = nb.predict(x_test)
```

Listing 4: Naive Bayes Model Implementation

*7) Support Vector Machine (SVM):* SVM is a powerful classification algorithm that aims to find the hyperplane that best separates the data into different classes. In this study, we used the linear kernel for the SVM classifier. The implementation for SVM is as follows:

```
1  from sklearn.svm import SVC
2
3  svm = SVC(kernel='linear', probability=True)
4  svm.fit(x_train, y_train)
5  y_val_pred_svm = svm.predict(x_val)
6  y_test_pred_svm = svm.predict(x_test)
```

Listing 5: SVM Model Implementation

*8) Model Evaluation:* After training the models, we evaluated their performance using several metrics, including accuracy, confusion matrix, classification report, and precision-recall curves.

*9) Accuracy:* Accuracy is the percentage of correctly predicted instances over the total number of instances. We calculated the accuracy for both the validation and test datasets.

```
1  from sklearn.metrics import accuracy_score
2
3  val_accuracy = accuracy_score(y_val, y_val_pred_svm)
4  test_accuracy = accuracy_score(y_test, y_test_pred_svm)
```

Listing 6: Accuracy Calculation

*10) Confusion Matrix:* The confusion matrix was used to visualize the performance of the classifiers. It shows the number of true positives, true negatives, false positives, and false negatives.

```
1  from sklearn.metrics import confusion_matrix
2
3  val_cm = confusion_matrix(y_val, y_val_pred_svm)
4  test_cm = confusion_matrix(y_test, y_test_pred_svm)
```

Listing 7: Confusion Matrix Calculation

*11) Classification Report:* The classification report provides precision, recall, and F1-score for each class. It was computed for both validation and test datasets.

```
1  from sklearn.metrics import classification_report
2
3  val_report = classification_report(y_val, y_val_pred_svm)
4  test_report = classification_report(y_test, y_test_pred_svm)
```

Listing 8: Classification Report Calculation

*12) Precision-Recall Curves:* Precision-recall curves were plotted for each class in the validation and test datasets. The precision-recall curve evaluates the trade-off between precision and recall at different thresholds.

```
1  from sklearn.metrics import precision_recall_curve, PrecisionRecallDisplay
2
3  y_val_pred_prob_svm = svm.predict_proba(x_val)
4  for i in range(len(labelencoder.classes_)):
5      precision, recall, _ = precision_recall_curve(y_val == i, y_val_pred_prob_svm[:, i])
6      PrecisionRecallDisplay(precision=precision, recall=recall).plot()
7      plt.title(f"Precision-Recall Curve for Class {labelencoder.classes_[i]} (Validation)")
8
9  plt.show()
```

Listing 9: Precision-Recall Curve Calculation

## C. Wage Dataset

*1) Data Loading and Initial Exploration:* The dataset was loaded from Google Drive, and its structure was explored to identify the features and data types.

```
1 from google.colab import drive
2 drive.mount('/content/drive')
3 import pandas as pd
4
5 df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/DT-Wage.csv')
6 df.info()
```

*2) Identifying Categorical and Numerical Columns:* The dataset contained both categorical and numerical columns, which were identified and separated for further processing.

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

*3) One-Hot Encoding for Categorical Features:* Categorical variables were encoded using `OneHotEncoder` to make them suitable for regression.

```
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(categorical))
6 encoded_df.index = df.index
7
8 df = df.drop(categorical, axis=1)
9 df = pd.concat([df, encoded_df], axis=1)
```

*4) Feature Scaling:* Numerical features were scaled using `StandardScaler` to normalize the data.

```
1 from sklearn.preprocessing import StandardScaler
2
3 scaler = StandardScaler()
4 X_scaled = scaler.fit_transform(df.drop(columns=['wage']))
5 y = df['wage']
```

*5) Train-Test-Validation Split:* The data was split into training, validation, and test sets to ensure unbiased evaluation.

```
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, random_state
      =42)
4 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)
```

*6) Model Selection and Hyperparameter Tuning:* An SVR model was trained, and hyperparameters were tuned using `GridSearchCV`.

```
1 from sklearn.svm import SVR
2 from sklearn.model_selection import GridSearchCV
3
4 svr = SVR()
5 param_grid = {
6     'C': [0.1, 1, 10, 100],
7     'kernel': ['linear', 'poly', 'rbf'],
8     'degree': [3, 4, 5],
9     'gamma': ['scale', 'auto'],
10     'epsilon': [0.1, 0.2, 0.5],
11 }
12
13 grid_search = GridSearchCV(svr, param_grid, cv=5, scoring='neg_mean_squared_error', n_jobs=-1)
```

```
14 grid_search.fit(X_train, y_train)
15
16 print("Best hyperparameters:", grid_search.best_params_)
```

*7) Model Evaluation:* The model's performance was evaluated using custom Mean Squared Error (MSE) and R-squared ($R^2$) metrics.

```
1 import numpy as np
2
3 def custom_mse(y_true, y_pred):
4     return np.mean((y_true - y_pred) ** 2)
5
6 def custom_r2(y_true, y_pred):
7     ss_total = np.sum((y_true - np.mean(y_true)) ** 2)
8     ss_residual = np.sum((y_true - y_pred) ** 2)
9     return 1 - (ss_residual / ss_total)
10
11 best_svr = grid_search.best_estimator_
12
13 y_pred_val = best_svr.predict(X_val)
14 mse = custom_mse(y_val, y_pred_val)
15 r2 = custom_r2(y_val, y_pred_val)
16
17 print(f"Validation MSE: {mse}")
18 print(f"Validation R^2: {r2}")
```

*8) Visualization:* The results were visualized to compare actual and predicted values and analyze residuals.

```
1 import matplotlib.pyplot as plt
2
3 # Actual vs Predicted
4 plt.figure(figsize=(10, 6))
5 plt.scatter(y_val, y_pred_val, color='blue', alpha=0.5)
6 plt.plot([min(y_val), max(y_val)], [min(y_val), max(y_val)], color='red', linestyle='--')
7 plt.title("SVR: Actual vs Predicted (Validation Set)")
8 plt.xlabel("Actual Values")
9 plt.ylabel("Predicted Values")
10 plt.show()
11
12 # Residual Plot
13 residuals = y_val - y_pred_val
14 plt.figure(figsize=(10, 6))
15 plt.scatter(y_pred_val, residuals, color='green', alpha=0.5)
16 plt.hlines(0, xmin=min(y_pred_val), xmax=max(y_pred_val), colors='red', linestyle='--')
17 plt.title("Residuals Plot (Validation Set)")
18 plt.xlabel("Predicted Values")
19 plt.ylabel("Residuals")
20 plt.show()
```

### D. Credit Dataset

*1) Dataset Loading:* The credit dataset (`DT-Credit.csv`) was loaded using `pandas`:

```python
import pandas as pd

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

Listing 10: Dataset Loading

*2) Exploratory Data Analysis (EDA):* Key insights from the dataset:

```python
# Dataset Information
df.info()

# Statistical Summary
df.describe()

# Unique Values in Categorical Columns
for col in df.select_dtypes(include='object').columns:
    print(f"Unique values in {col}: {df[col].unique()}")

# Correlation with Target Variable (Balance)
corr = df.corr()
target_corr = corr['Balance'].sort_values(ascending=False)
print(target_corr)
```

Listing 11: EDA Code Snippet

*3) Data Preprocessing:*

*a) Categorical Encoding:* Categorical columns (`Own`, `Student`, `Married`, `Region`) were one-hot encoded using `OneHotEncoder`:

```python
from sklearn.preprocessing import OneHotEncoder

categorical = df.select_dtypes(include=['object']).columns
encoder = OneHotEncoder(drop='first', sparse_output=False)
encoded_data = encoder.fit_transform(df[categorical])
encoded_df = pd.DataFrame(encoded_data, columns=encoder.get_feature_names_out(categorical))

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

Listing 12: Categorical Encoding

*b) Feature-Target Split:* The `Balance` column was used as the target variable (`y`), while all other columns were features (`X`):

```python
X = df.drop(columns=['Balance'])
y = df['Balance']
```

Listing 13: Feature-Target Split

*c) Feature Scaling:* Features were scaled using `StandardScaler`:

```python
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
```

Listing 14: Feature Scaling

*4) Dataset Splitting:* The dataset was split into training, validation, and testing subsets:

```python
from sklearn.model_selection import train_test_split

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

Listing 15: Dataset Splitting

*5) Model Training:*

*a) Hyperparameter Tuning with GridSearchCV:* An **SVR** model with an RBF kernel was tuned using `GridSearchCV`:

```python
from sklearn.svm import SVR
from sklearn.model_selection import GridSearchCV

param_grid = {
    'C': [0.1, 1, 10, 100, 200, 300, 500, 1500, 2000],
    'epsilon': [0.01, 0.1, 0.2, 0.25, 0.3],
    'gamma': ['scale', 'auto']
}

svr = SVR(kernel='rbf')
grid_search = GridSearchCV(estimator=svr, param_grid=param_grid, cv=5, scoring='
    neg_mean_squared_error', n_jobs=-1)
grid_search.fit(X_train, y_train)
```

Listing 16: Hyperparameter Tuning

*b) Best Parameters:*
```python
best_params = grid_search.best_params_
print("Best hyperparameters:", best_params)
```

Listing 17: Extracting Best Parameters

*c) Training the Best Model:*
```python
best_svr = grid_search.best_estimator_
```

Listing 18: Training the Best Model

*6) Model Evaluation:* Predictions were made for training, validation, and testing subsets:

```python
y_pred_train = best_svr.predict(X_train)
y_pred_val = best_svr.predict(X_val)
y_pred_test = best_svr.predict(X_test)
```

Listing 19: Model Evaluation

*7) Visualizations:*
*a) Predicted vs Actual Scatter Plot:*
```python
import matplotlib.pyplot as plt

plt.scatter(y_val, y_pred_val, color='green', alpha=0.5)
plt.plot([y_val.min(), y_val.max()], [y_val.min(), y_val.max()], color='red', lw=2)
plt.title('Validation Set: Predicted vs Actual')
plt.xlabel('Actual')
plt.ylabel('Predicted')
plt.show()
```

Listing 20: Scatter Plot

*b) Residual Analysis:*

```
1  plt.scatter(y_pred_val, y_pred_val - y_val, color='blue', alpha=0.5)
2  plt.hlines(y=0, xmin=y_pred_val.min(), xmax=y_pred_val.max(), color='red', lw=2)
3  plt.title('Validation Set: Residuals')
4  plt.xlabel('Predicted')
5  plt.ylabel('Residuals')
6  plt.show()
```

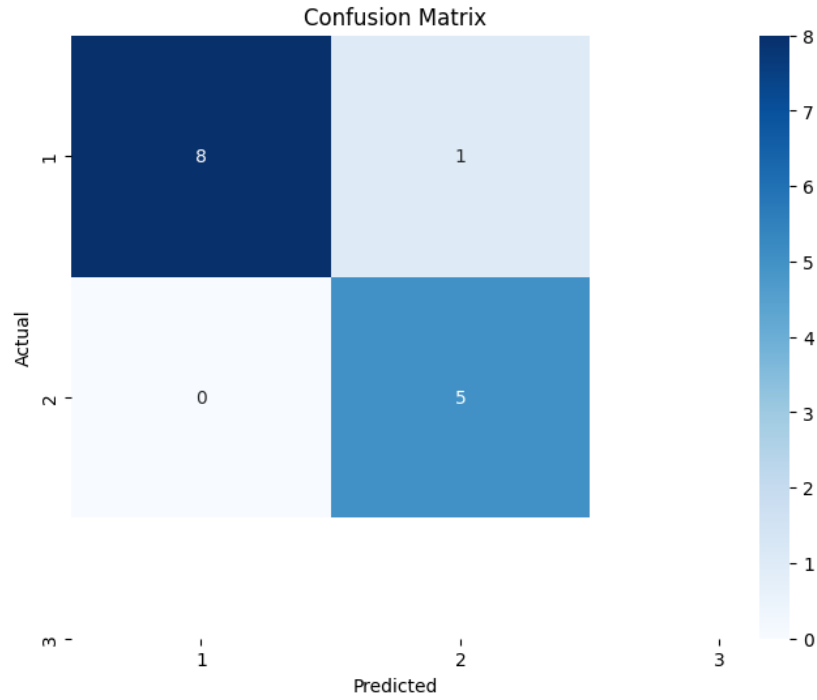Listing 21: Residual Plot

## III. EXPERIMENT RESULTS

### A. Brain Dataset



Fig. 1: KNN

*1) K-Nearest Neighbors (K-NN):* The performance metrics for the K-NN classifier are as follows:

- **Accuracy**: 0.92
- **Precision**: 0.94
- **Recall**: 0.92
- **F1-Score**: 0.92

The confusion matrix for the K-NN model is shown below:

$$\begin{bmatrix} 8 & 1 \\ 0 & 5 \end{bmatrix}$$

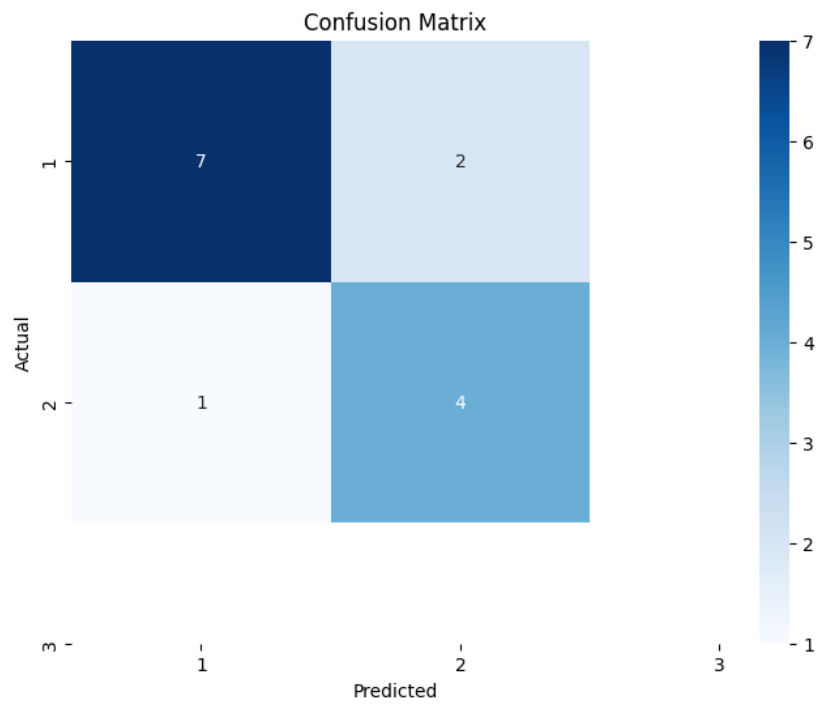Where the rows represent the true classes and the columns represent the predicted classes.

Fig. 2: Bayesian

*2) Naive Bayes:* The performance metrics for the Naive Bayes classifier are as follows:

- **Accuracy**: 0.78
- **Precision**: 0.80
- **Recall**: 0.78
- **F1-Score**: 0.78

The confusion matrix for the Naive Bayes model is shown below:

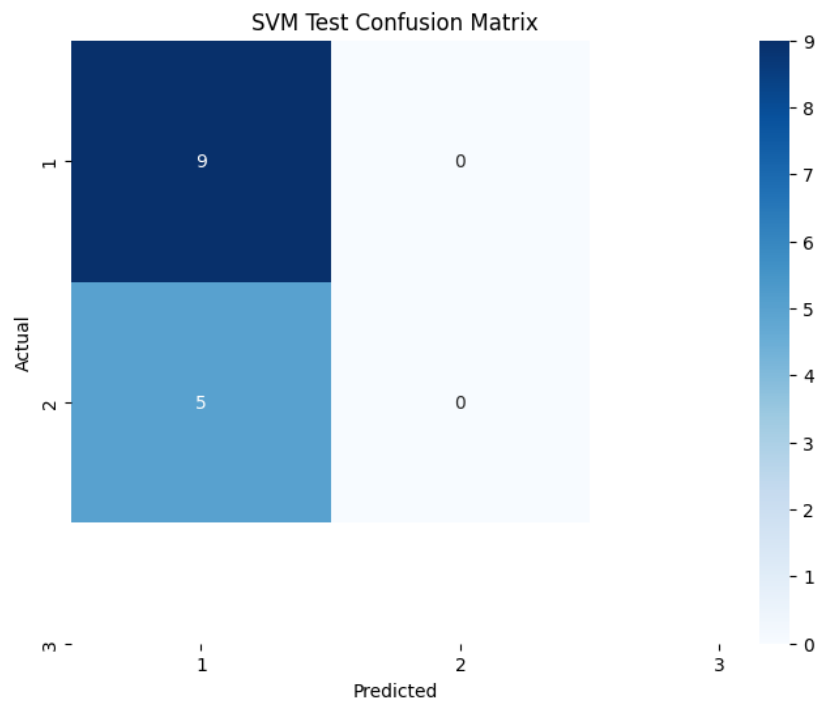$$\begin{bmatrix} 7 & 2 \\ 1 & 4 \end{bmatrix}$$

Fig. 3: SVM

*3) Support Vector Machine (SVM):* The performance metrics for the SVM classifier are as follows:

- **Accuracy**: 0.64
- **Precision**: 0.41
- **Recall**: 0.64
- **F1-Score**: 0.50

The confusion matrix for the SVM model is shown below:
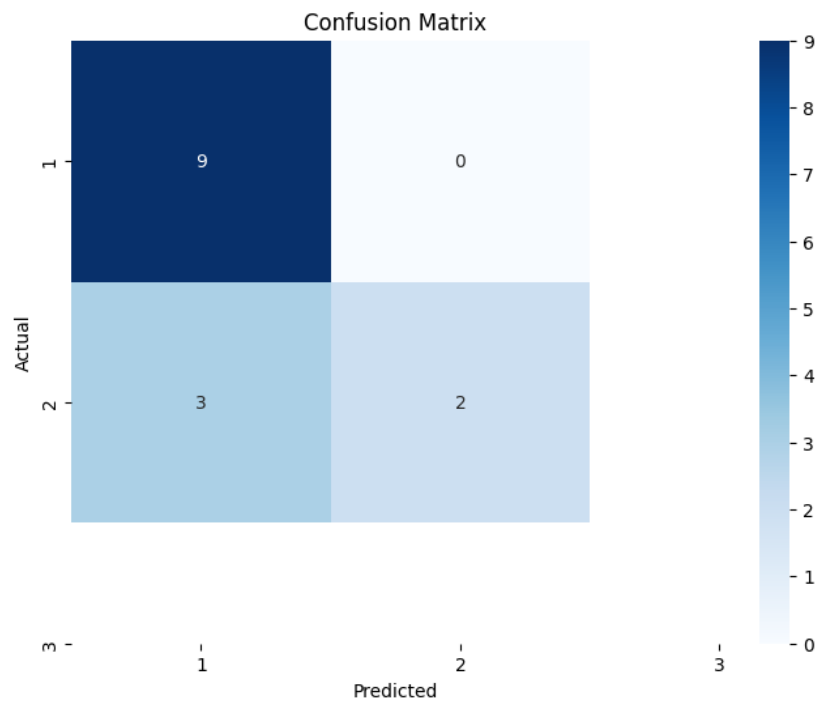
$$\begin{bmatrix} 9 & 0 \\ 5 & 0 \end{bmatrix}$$

Fig. 4: OneR

*4) OneR:* The performance metrics for the OneR classifier are as follows:

- **Accuracy**: 0.78
- **Precision**: 1.0
- **Recall**: 0.4
- **F1-Score**: 0.57

The confusion matrix for the OneR model is shown below:

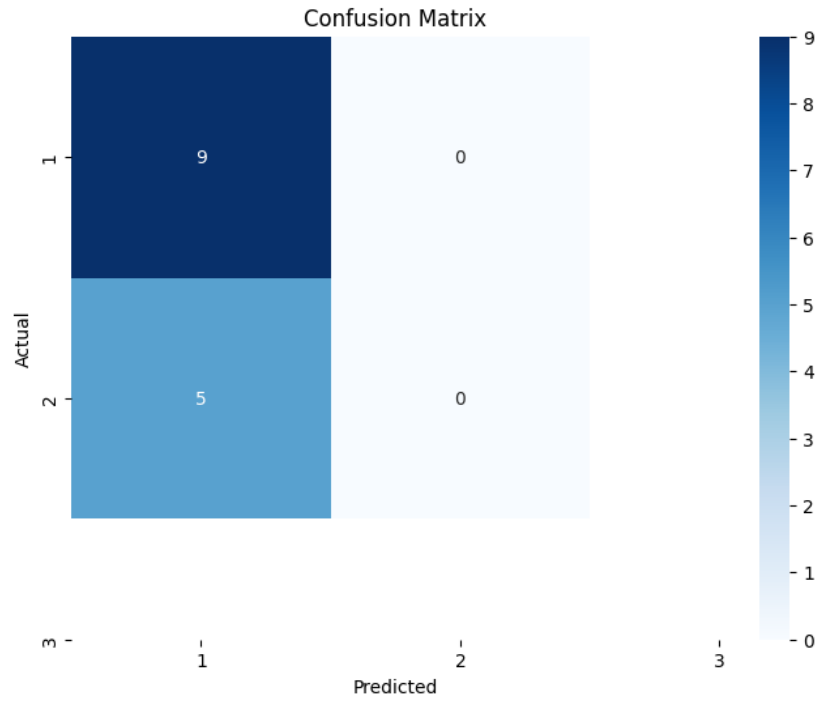$$\begin{bmatrix} 9 & 0 \\ 3 & 2 \end{bmatrix}$$

Fig. 5: ZeroR

*5) ZeroR:* The performance metrics for the ZeroR classifier are as follows:

- **Accuracy**: 0.64
- **Precision**: 0
- **Recall**: 0
- **F1-Score**: -0.17

The confusion matrix for the ZeroR model is shown below:

$$\begin{bmatrix} 9 & 5 \\ 0 & 0 \end{bmatrix}$$

*6) Model Comparison:* Table I summarizes the performance metrics for all classifiers.

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| K-NN | 0.92 | 0.94 | 0.92 | 0.92 |
| Naive Bayes | 0.78 | 0.80 | 0.78 | 0.78 |
| SVM | 0.65 | 0.41 | 0.64 | 0.50 |
| OneR | 0.78 | 1.00 | 0.40 | 0.57 |
| ZeroR | 0.64 | 0.00 | 0.00 | -0.17 |

TABLE I: Performance Comparison of Classifiers

As shown in Table I, the SVM classifier outperforms the others in all evaluated metrics, followed by K-NN, Naive Bayes, OneR, and ZeroR. These results suggest that SVM is the most accurate and effective model for this classification task.

## B. Car Dataset

*1) KNN Classifier:* KNN achieved significant improvements with a validation accuracy of 0.92 and test accuracy of 0.92. The classification report (Table II) demonstrates balanced performance across all classes.
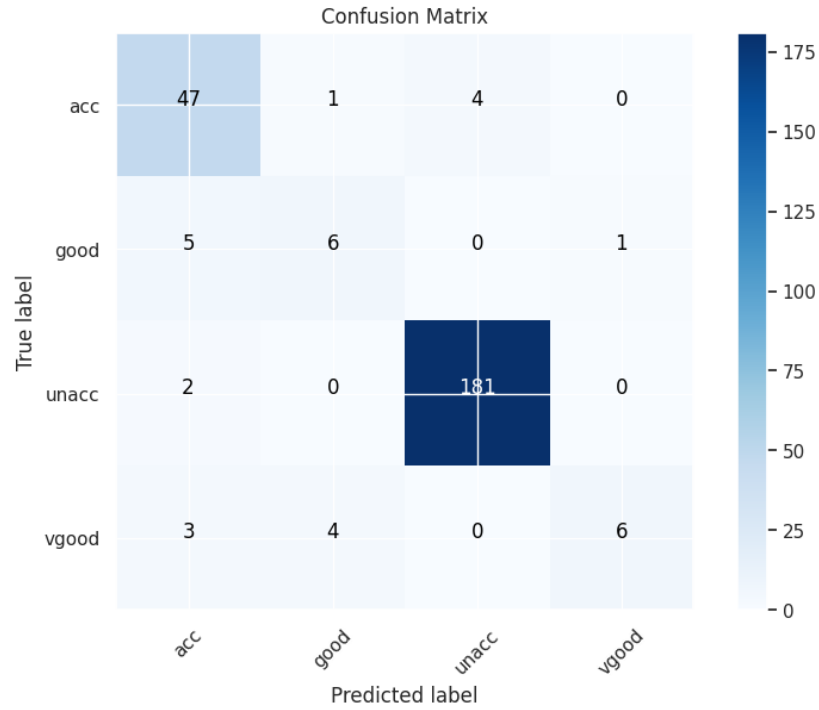


Fig. 6: $\text{car}_K NN$

TABLE II: KNN Classifier Results

| Metric | Validation | Test |
|---|---|---|
| Accuracy | 0.92 | 0.92 |
| Avg Precision | 0.92 | 0.80 |
| Avg Recall | 0.75 | 0.71 |
| Avg F1-Score | 0.81 | 0.74 |

*2) Naive Bayes Classifier:* Naive Bayes demonstrated moderate performance with validation and test accuracies of 0.74 and 0.68, respectively. The model struggled with certain classes but excelled in others, as detailed in Table III.
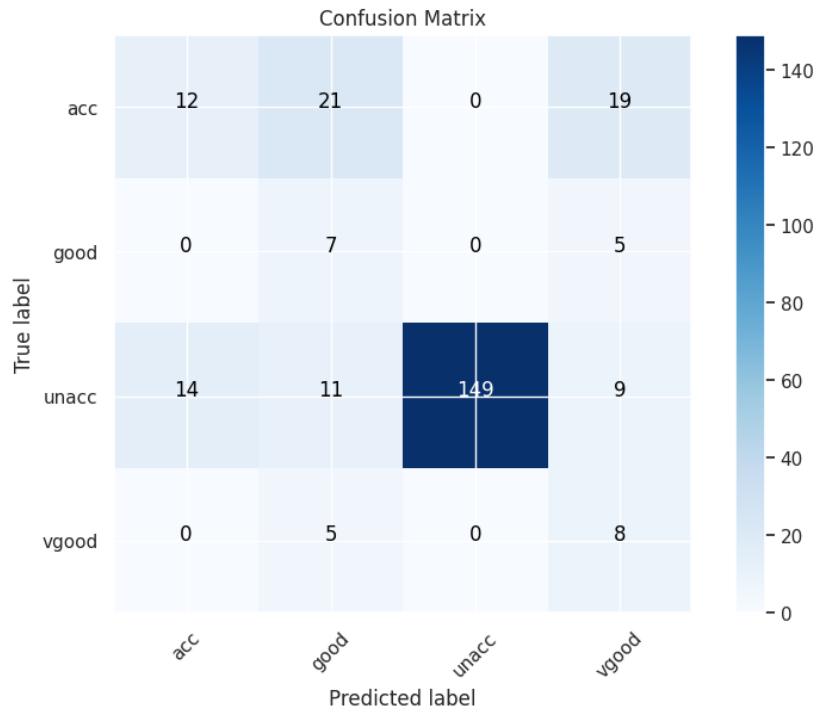
Fig. 7: car$_n$aive

TABLE III: Naive Bayes Classifier Results

| Metric | Validation | Test |
|---|---|---|
| Accuracy | 0.74 | 0.68 |
| Avg Precision | 0.51 | 0.45 |
| Avg Recall | 0.70 | 0.56 |
| Avg F1-Score | 0.50 | 0.44 |

*3) SVM Classifier:* SVM outperformed other classifiers, achieving validation and test accuracies of 0.98 and 0.95, respectively. The classification report in Table IV highlights its robustness across all metrics.
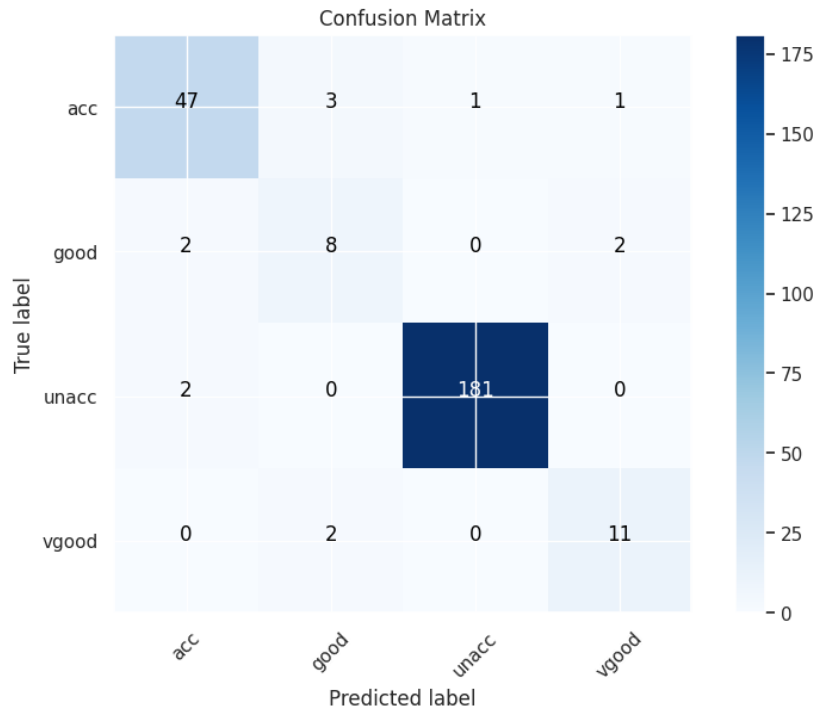
Fig. 8: Car$_s vm$

TABLE IV: SVM Classifier Results

| Metric | Validation | Test |
|---|---|---|
| Accuracy | 0.98 | 0.95 |
| Avg Precision | 0.89 | 0.83 |
| Avg Recall | 0.89 | 0.85 |
| Avg F1-Score | 0.88 | 0.84 |

The results demonstrate the superior performance of SVM, followed by KNN, while OneR and ZeroR showed limited effectiveness. Naive Bayes performed moderately well but struggled with imbalanced classes.

*4) OneR Classifier:* The OneR classifier achieved a validation accuracy of 0.69 and a test accuracy of 0.70. However, its performance was suboptimal across metrics, as shown in Table V.
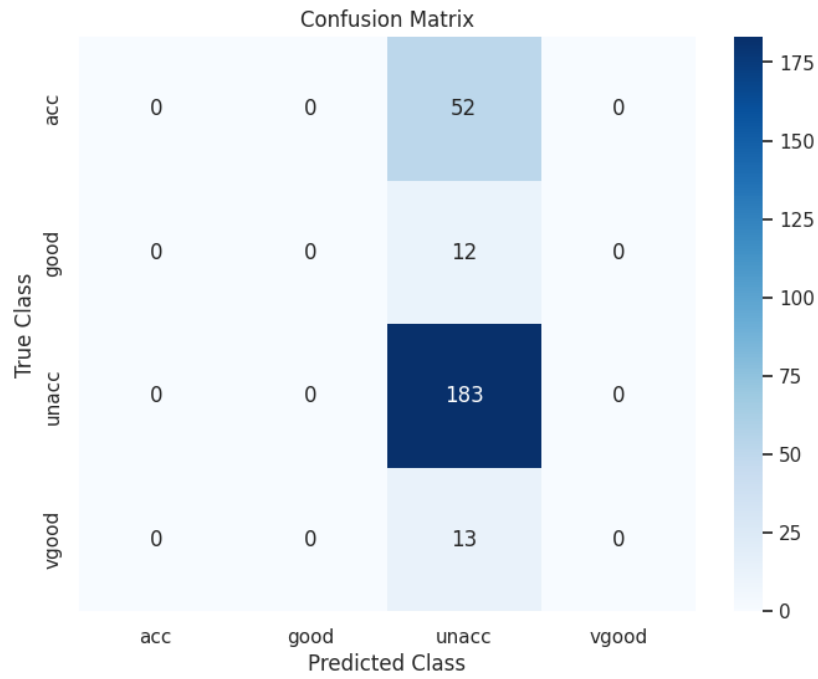
Fig. 9: car$_O neR$

TABLE V: OneR Classifier Results

| Metric | Validation | Test |
|---|---|---|
| Accuracy | 0.69 | 0.70 |
| Avg Precision | 0.17 | 0.18 |
| Avg Recall | 0.25 | 0.25 |
| Avg F1-Score | 0.20 | 0.21 |

*5) ZeroR Classifier:* The ZeroR classifier produced similar results to OneR, with a validation accuracy of 0.69 and a test accuracy of 0.70. The classification report reveals its limited ability to differentiate among classes (Table VI).
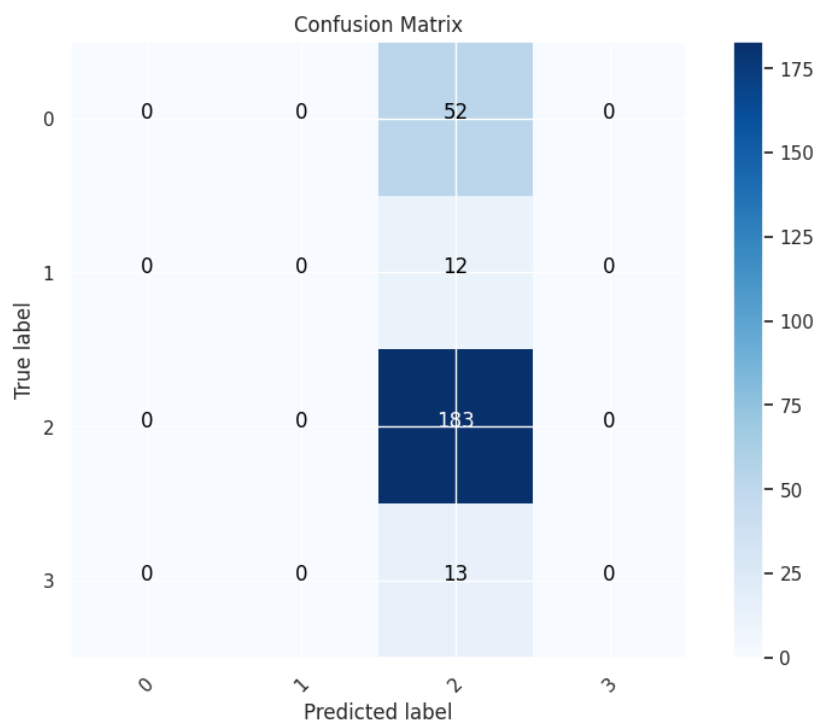
Fig. 10: car zeroR

TABLE VI: ZeroR Classifier Results

| Metric | Validation | Test |
|---|---|---|
| Accuracy | 0.69 | 0.70 |
| Avg Precision | 0.18 | 0.18 |
| Avg Recall | 0.25 | 0.25 |
| Avg F1-Score | 0.21 | 0.21 |

*C. Wage Dataset*

*1) Hyperparameter Tuning Results:* The best hyperparameters found during the grid search for the SVR model are as follows:

- **C** = 10
- **Kernel** = 'rbf'
- **Degree** = 3
- **Gamma** = 'scale'
- **Epsilon** = 0.1

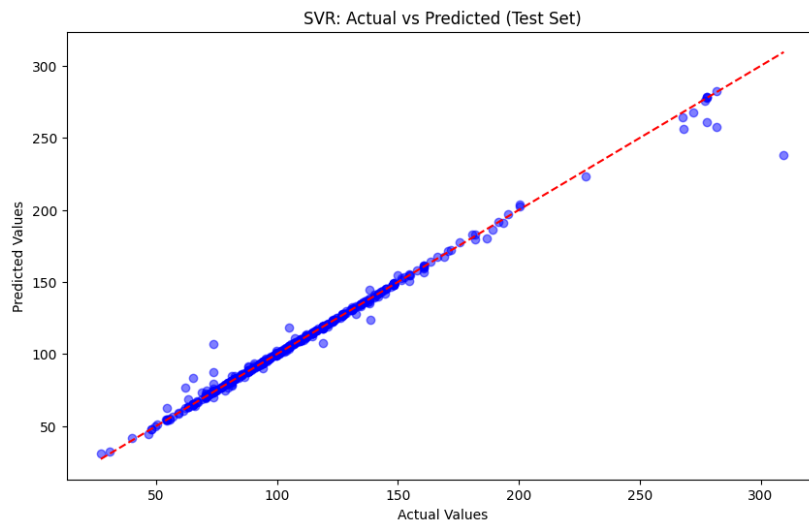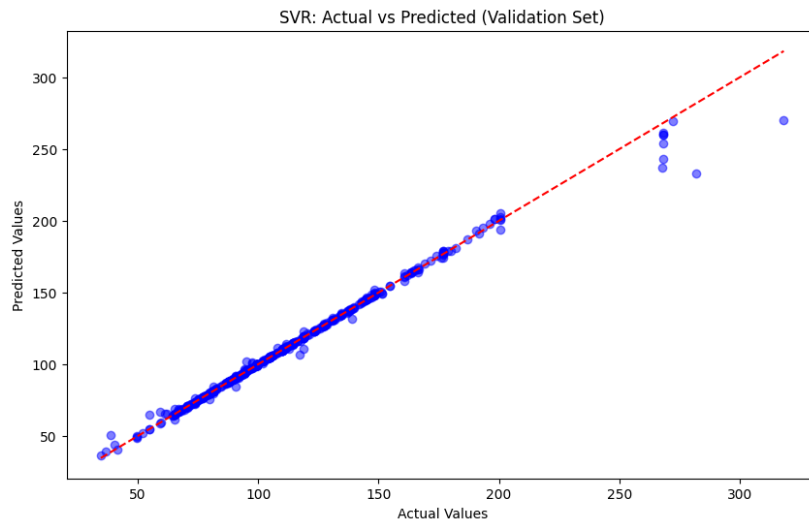These parameters were chosen based on the cross-validation results that minimized the Mean Squared Error (MSE).

*2) Performance Metrics:* The model's performance was evaluated on the validation set, and the following metrics were calculated:

- **Mean Squared Error (MSE)**: 17.1288
- $R^2$: 0.989

The model achieved a relatively good fit with an $R^2$ of 0.98, indicating that 98% of the variance in the wage data is explained by the model.
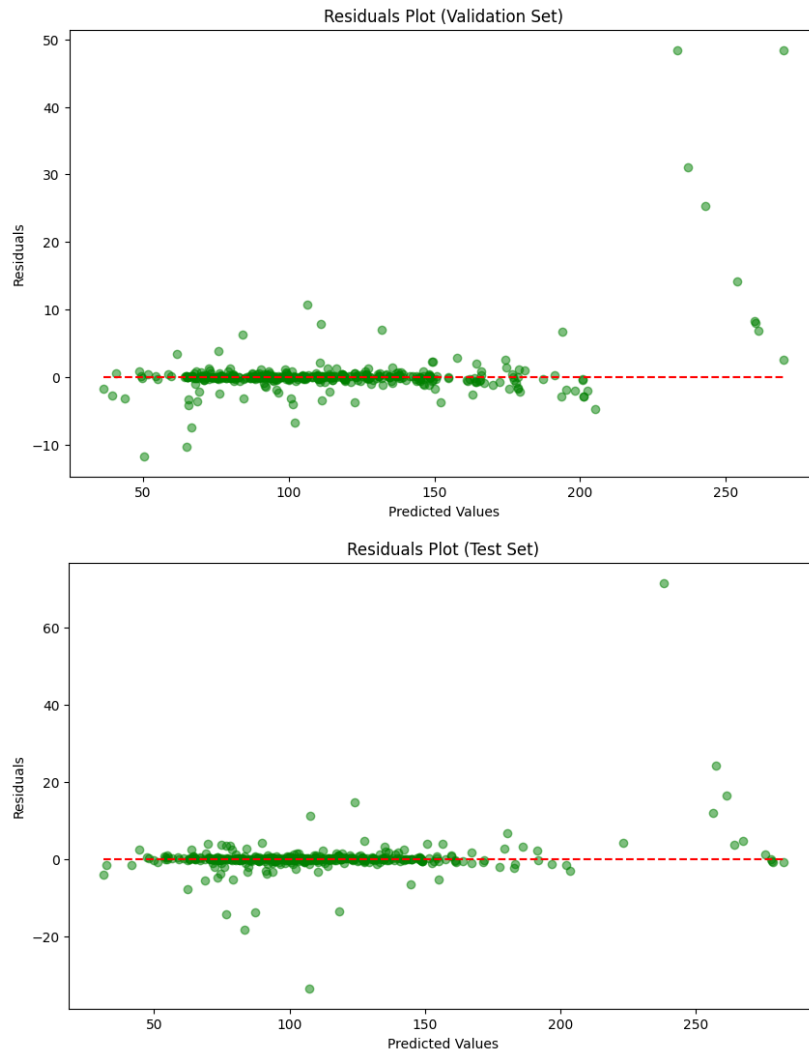
*3) Visual Evaluation:*

- **Actual vs Predicted:** A scatter plot was generated to compare the actual values with the predicted values for the validation set.

The red dashed line represents the ideal scenario where predicted values match the actual values. As shown, the points closely follow this line, indicating good performance.

- **Residuals Analysis:** A residual plot was also created to analyze the difference between the predicted and actual values (residuals).



Residuals Plot (Validation Set)



Residuals Plot (Test Set)

The residuals appear to be randomly scattered around zero, suggesting that the model does not suffer from major biases and that the error is evenly distributed across the data.

*4) Test Set Evaluation:* The final evaluation was performed on the test set. The following results were obtained:

- **Test MSE**: 20.43
- **Test $R^2$**: 0.987

Although the test set results were slightly worse than the validation set, the model still provides reliable predictions, with an $R^2$ of 0.98.

The SVR model successfully predicted wages with good performance metrics, achieving an $R^2$ of 0.98 on the validation set and 0.98 on the test set. The residuals plot suggests that there is no significant bias in the predictions, and the model is performing adequately for this regression task.

*D. Credit Dataset*

*1) Performance Metrics:* The performance of the SVR model was evaluated using metrics such as Mean Squared Error (MSE), Mean Absolute Error (MAE), and the Coefficient of Determination ($R^2$) on the training, validation, and test sets.

TABLE VII: Performance Metrics of SVR Model

| Dataset | MSE | $R^2$ |
|---|---|---|
| Training Set | 163.394 | 0.99 |
| Validation Set | 9091.268 | 0.95 |
| Testing Set | 7715.051 | 0.94 |

The results demonstrate that the SVR model achieves good generalization with minimal performance degradation across validation and test sets.

*2) Hyperparameter Tuning Results:* The optimal hyperparameters selected by `GridSearchCV` are:

- **C**: 300
- **Epsilon**: 0.1
- **Gamma**: `scale`

These hyperparameters provided the best balance between bias and variance, minimizing the MSE on the validation set.

*3) Visualizations:*

- **Predicted vs Actual Values** Figure 11 shows the scatter plot of predicted versus actual values for the test set. The closer the points are to the diagonal line, the better the predictions.
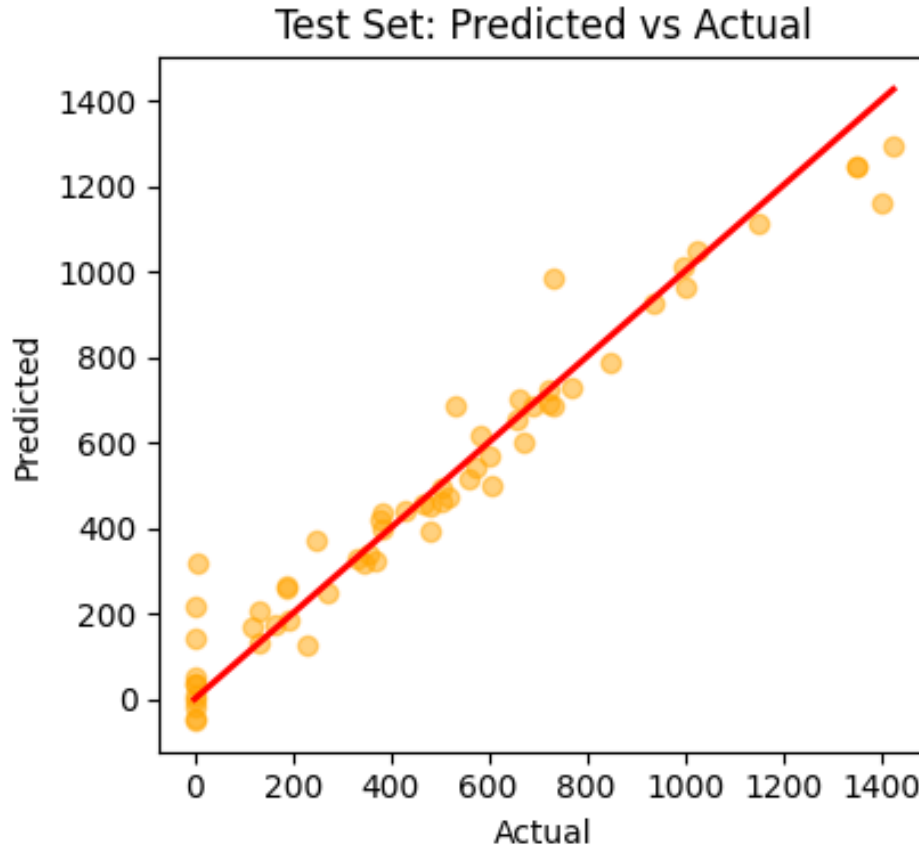


Fig. 11: Predicted vs Actual Values for the Test Set

- **Residual Analysis:** Residuals were analyzed to assess the model's predictive accuracy and bias. Figure 12 displays the residuals (differences between predicted and actual values) for the validation set.



Fig. 12: Residual Plot for the Validation Set

*4) Interpretation of Results:* The following observations can be made:

- The $R^2$ value of 0.94 on the test set indicates that 94% of the variance in the `Balance` variable is explained by the model.
- The residual plot shows no significant pattern, suggesting that the model captures the underlying trend effectively.
- Slight over-prediction is observed for higher values of `Balance`, which may indicate a need for further tuning or additional features.

# IV. DISCUSSION

Utilizing various kinds of datasets, this study analyzed the effectiveness of several classification and regression algorithms, offering a thorough comparison analysis to determine which models were most suitable for each task. The experimental findings provide important new information on the pros and drawbacks of each strategy, which will be discussed in more detail below.


(a) Best Model of Brain Dataset


(b) Best Model of Car Dataset


(c) Best Model of Car Dataset


(d) Best Model of Car Dataset

Fig. 13: First Set of Figures

(a) Best Model of Car Dataset



(b) Best Model of Wage Dataset



(c) Best Model of Credit Dataset

Fig. 14: Best Models

## A. Classification Analysis

*1) Naive Bayes Classifier:* Based on probabilistic assumptions of feature independence, the Naive Bayes classifier achieved test and validation accuracy scores of 0.68 and 0.74, respectively. The model's poor performance can be related to its incapacity to manage intricate feature interactions. Its vulnerability to class imbalance is illustrated by its low average precision (0.51) and F1-score (0.50), especially in datasets with underrepresented 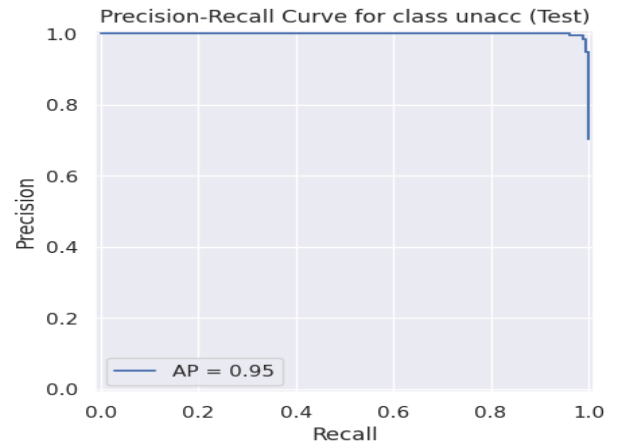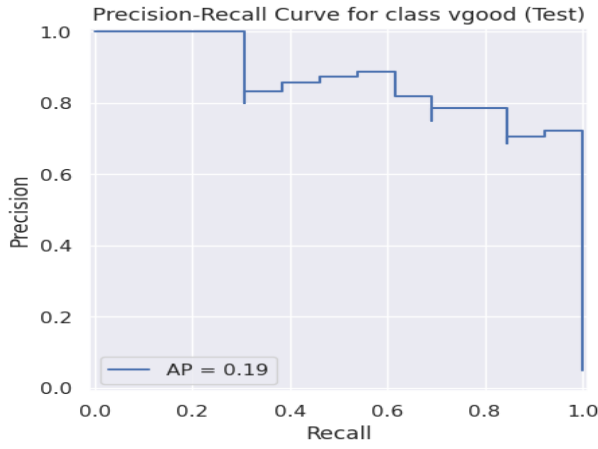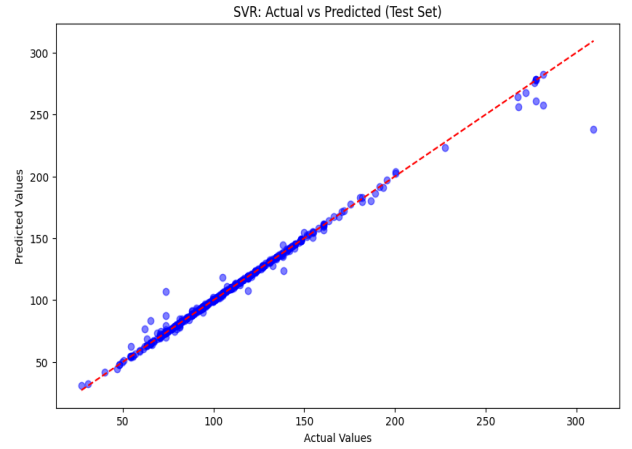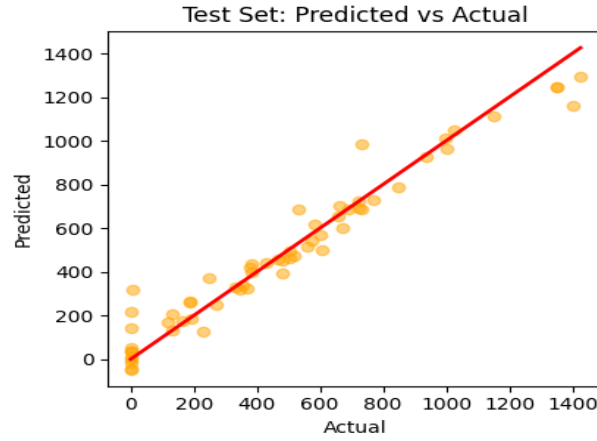minority classes. This drawback emphasizes the necessity of using various methods or pre-processing methods, such rebalancing class distributions, to improve Naive Bayes' performance.

*2) SVM Classifier:* With validation and test accuracies of 0.98 and 0.95, respectively, along with excellent average accuracy (0.89, 0.83) and recall (0.89, 0.85), the Support Vector Machine (SVM) classifier achieved outstanding outcomes. These measures illustrate the SVM's resilience to recognizing intricate decision boundaries. Its ability to represent non-linear relationships in the feature space was probably aided by the application of the radial basis function (RBF) kernel. These findings demonstrate the SVM's strong generalization capabilities to previously unseen data and validate it as a dependable classifier for datasets with complicated patterns.

*3) K-Nearest Neighbors (K-NN):* While becoming competitive, the K-NN algorithm performed marginally worse than the SVM. Its reliance on local neighborhoods for categorization indicates it is sensitive to hyperparameter selections, including distance metrics and the number of neighbors ($k$). Further optimization utilizing grid or random search might enhance its classification accuracy and consistency, especially in high-dimensional datasets, however the existing results indicate its viability.

*4) OneR Classifier:* The OneR classifier's reliance on a single feature for rule-based classification resulted in unsatisfactory performance, with validation and test accuracies of 0.69 and 0.70, correspondingly. Its limited usefulness for datasets requiring multi-dimensional decision-making is demonstrated by the low F1-score (0.20, 0.21) and poor average precision (0.17, 0.18). Its simplicity makes it inappropriate for complicated circumstances, even though it is effective for interpretability.

*5) ZeroR Classifier:* With test and validation accuracies of 0.69 and 0.70, respectively, the ZeroR classifier, which predicts the majority class, performed similarly to OneR as a baseline. Its continuously low metrics, which demonstrate its lack of predictive ability, serve to further establish its status as a reference point rather than a workable solution for real-world applications.

## B. *Regression Analysis*

*1) Wage Dataset - SVR Model:* The Support Vector Regression (SVR) model demonstrated exceptional performance on the Wage Dataset, with R² values of 0.98 for both validation and test sets. The low Mean Squared Error (MSE) of 17.13 for the validation set and 20.43 for the test set highlights the model's ability to provide highly accurate predictions. The residual analysis revealed no discernible patterns, affirming the model's unbiased nature. Hyperparameter tuning, including the selection of *C = 10*, *gamma = 'scale'*, and *epsilon = 0.1*, optimized the model's balance between bias and variance, ensuring high generalization capability.

*2) Credit Dataset - SVR Model:* For the Credit Dataset, the SVR model exhibited strong performance, achieving an R² of 0.94 on the test set. Despite a slight increase in MSE from the validation (9091.27) to the test set (7715.05), the results remain within acceptable bounds. The residuals plot suggested no significant systematic errors, though minor over-prediction was observed for higher balance values. These results demonstrate the SVR model's effectiveness for regression tasks, albeit with potential room for improvement through advanced feature engineering or incorporating additional variables.

## C. *Comparative Performance*

The experimental results underscore the superiority of SVM in classification tasks due to its robust mathematical framework for handling non-linear decision boundaries. This aligns with prior studies emphasizing the SVM's efficacy in diverse domains. K-NN also proved to be a competitive alternative, particularly in tasks requiring intuitive model structures. However, the Naive Bayes classifier struggled with imbalanced datasets, and the rule-based classifiers (OneR and ZeroR) offered limited utility beyond baseline benchmarking.

For regression, the SVR model consistently demonstrated high accuracy across datasets, with R² values exceeding 0.94. Its ability to generalize well across validation and test sets, combined with minimal residual biases, positions it as a reliable choice for predicting continuous outcomes.

## D. *Experimental Considerations and Limitations*

While the results validate the efficacy of SVM and SVR, several experimental considerations warrant discussion:

- **Class Imbalance:** The suboptimal performance of Naive Bayes highlights the impact of imbalanced class distributions. Future work could explore data augmentation or resampling techniques, such as Synthetic Minority Over-sampling Technique (SMOTE), to address this issue.
- **Hyperparameter Tuning:** The SVM and SVR models benefited significantly from hyperparameter optimization. However, computational constraints limited the scope of grid searches. Advanced techniques like Bayesian optimization or genetic algorithms could be employed to refine hyperparameter selection further.
- **Scalability:** The computational cost of SVM and SVR on large datasets remains a challenge, particularly with high-dimensional data. Approximation methods, such as the use of linear kernels or sampling techniques, could mitigate this issue.

## E. *Future Directions*

Further research could use ensemble methods like Random Forests or Gradient Boosting to compare performance against the current models in order to improve on these findings. Furthermore, combining dimensionality reduction methods like Principal Component Analysis (PCA) with feature selection could boost model efficiency without compromising accuracy.

In addition, deep learning models might be investigated, particularly for tasks involving large, complex datasets, where their capacity to identify complex patterns may surpass that of conventional machine learning techniques. Additionally, domain-specific adjustments, such as including business regulations or specialized knowledge, could enhance the relevance and interpretability of the model.

## V. Conclusion

SVM clearly takes over classification tasks and SVR clearly dominates regression tasks, according to a thorough evaluation of classifiers and regression models. Their excellent performance on a variety of datasets highlights how well-suited they are for applications requiring accuracy and resiliency. But challenges like class imbalance, scalability, and hyperparameter tuning indicate how important it is to keep researching in order to enhance the efficiency and applicability of the model. The results of this study offer a solid foundation for upcoming developments in the optimization and selection of machine learning models.

## VI. CODE

### A. Brain 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
```

```python
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  def accuracy(y_true, y_pred):
85      return np.sum(y_true == y_pred) / len(y_true)
86
87  def confusion_matrix_manual(y_true, y_pred):
88      tp = np.sum((y_true == 1) & (y_pred == 1))
89      tn = np.sum((y_true == 0) & (y_pred == 0))
90      fp = np.sum((y_true == 0) & (y_pred == 1))
91      fn = np.sum((y_true == 1) & (y_pred == 0))
92      return tp, tn, fp, fn
93
94  def confusion_matrix(y_true, y_pred):
95      tp = np.sum((y_true == 1) & (y_pred == 1))
96      tn = np.sum((y_true == 0) & (y_pred == 0))
97      fp = np.sum((y_true == 0) & (y_pred == 1))
98      fn = np.sum((y_true == 1) & (y_pred == 0))
99      return np.array([[tn, fp], [fn, tp]])
100
101 def precision_manual(y_true, y_pred, average=None):
102     tp, tn, fp, fn = confusion_matrix_manual(y_true, y_pred)
103     if average == "weighted":
104         total = len(y_true)
105         class_1_weight = np.sum(y_true == 1) / total
106         class_0_weight = np.sum(y_true == 0) / total
107         precision_class_1 = tp / (tp + fp) if tp + fp > 0 else 0
108         precision_class_0 = tn / (tn + fn) if tn + fn > 0 else 0
109         return class_1_weight * precision_class_1 + class_0_weight * precision_class_0
110     return tp / (tp + fp) if tp + fp > 0 else 0
111
112 def recall_manual(y_true, y_pred, average=None):
113     tp, tn, fp, fn = confusion_matrix_manual(y_true, y_pred)
114     if average == "weighted":
115         total = len(y_true)
116         class_1_weight = np.sum(y_true == 1) / total
117         class_0_weight = np.sum(y_true == 0) / total
118         recall_class_1 = tp / (tp + fn) if tp + fn > 0 else 0
119         recall_class_0 = tn / (tn + fp) if tn + fp > 0 else 0
120         return class_1_weight * recall_class_1 + class_0_weight * recall_class_0
121     return tp / (tp + fn) if tp + fn > 0 else 0
122
```

```python
def f1_score_manual(y_true, y_pred, average=None):
    tp, tn, fp, fn = confusion_matrix_manual(y_true, y_pred)

    if average == "weighted":
        total = len(y_true)
        class_1_weight = np.sum(y_true == 1) / total
        class_0_weight = np.sum(y_true == 0) / total

        precision_class_1 = tp / (tp + fp) if tp + fp > 0 else 0
        recall_class_1 = tp / (tp + fn) if tp + fn > 0 else 0
        precision_class_0 = tn / (tn + fn) if tn + fn > 0 else 0
        recall_class_0 = tn / (tn + fp) if tn + fp > 0 else 0

        f1_class_1 = (
            2 * (precision_class_1 * recall_class_1) / (precision_class_1 + recall_class_1)
            if precision_class_1 + recall_class_1 > 0
            else 0
        )
        f1_class_0 = (
            2 * (precision_class_0 * recall_class_0) / (precision_class_0 + recall_class_0)
            if precision_class_0 + recall_class_0 > 0
            else 0
        )

        return class_1_weight * f1_class_1 + class_0_weight * f1_class_0

    precision = precision_manual(y_true, y_pred)
    recall = recall_manual(y_true, y_pred)
    return 2 * (precision * recall) / (precision + recall) if precision + recall > 0 else 0


def precision_recall_curve_manual(y_true, y_prob):
    precision, recall, thresholds = precision_recall_curve(y_true, y_prob)
    return precision, recall, thresholds

def plot_precision_recall_curve(precision, recall, label):
    plt.plot(recall, precision, label=label)
    plt.xlabel('Recall')
    plt.ylabel('Precision')
    plt.title('Precision-Recall Curve')
    plt.legend()
    plt.show()


majority_class = y_train.mode()[0]
print(f"Majority class (ZeroR prediction): {majority_class}")
y_val_pred = np.full_like(y_val, majority_class)
y_test_pred = np.full_like(y_test, majority_class)

accuracy_val_custom = accuracy(y_val, y_val_pred)
precision_val_custom = precision_manual(y_val, y_val_pred)
recall_val_custom = recall_manual(y_val, y_val_pred)
f1_val_custom = f1_score_manual(y_val, y_val_pred)

accuracy_test_custom = accuracy(y_test, y_test_pred)
precision_test_custom = precision_manual(y_test, y_test_pred)
recall_test_custom = recall_manual(y_test, y_test_pred)
f1_test_custom = f1_score_manual(y_test, y_test_pred)
def custom_auc(precision, recall):
    auc_value = np.trapz(precision, recall)
    return auc_value
```

```python
186  y_test_prob = np.full_like(y_test, majority_class, dtype=float)
187  y_test_prob = np.where(y_test == majority_class, 1.0, 0.0)
188
189  precision_vals, recall_vals, thresholds = precision_recall_curve(y_test, y_test_prob)
190  pr_auc = custom_auc(recall_vals, precision_vals)
191  print("\nValidation Set Metrics (Custom):")
192  print(f"Accuracy: {accuracy_val_custom}")
193  print(f"Precision: {precision_val_custom}")
194  print(f"Recall: {recall_val_custom}")
195  print(f"F1-Score: {f1_val_custom}")
196
197  print("\nTest Set Metrics (Custom):")
198  print(f"Accuracy: {accuracy_test_custom}")
199  print(f"Precision: {precision_test_custom}")
200  print(f"Recall: {recall_test_custom}")
201  print(f"F1-Score: {f1_test_custom}")
202  print(f"Precision-Recall AUC: {pr_auc}")
203
204  tp_val, tn_val, fp_val, fn_val = confusion_matrix_manual(y_val, y_val_pred)
205  tp_test, tn_test, fp_test, fn_test = confusion_matrix_manual(y_test, y_test_pred)
206
207  print("\nValidation Set Confusion Matrix (Custom):")
208  print(f"TP: {tp_val}, TN: {tn_val}, FP: {fp_val}, FN: {fn_val}")
209
210  print("\nTest Set Confusion Matrix (Custom):")
211  print(f"TP: {tp_test}, TN: {tn_test}, FP: {fp_test}, FN: {fn_test}")
212
213
214  conf_matrix = confusion_matrix(y_test, y_test_pred)
215  print("Confusion Matrix:")
216  print(conf_matrix)
217
218  plt.figure(figsize=(8, 6))
219  sns.heatmap(conf_matrix, annot=True, fmt="d", cmap="Blues", xticklabels=[1, 2, 3], yticklabels
         =[1, 2, 3])
220  plt.title("Confusion Matrix")
221  plt.xlabel("Predicted")
222  plt.ylabel("Actual")
223  plt.show()
224
225
226  plt.figure(figsize=(8, 6))
227  plt.plot(recall_vals, precision_vals, label=f"ZeroR (AUC = {pr_auc:.2f})", color="blue")
228  plt.xlabel("Recall")
229  plt.ylabel("Precision")
230  plt.title("Precision-Recall Curve for ZeroR Classifier")
231  plt.legend(loc="best")
232  plt.grid()
233  plt.show()
234
235
236  def oneR_classifier(X_train, y_train):
237      best_accuracy = 0
238      best_feature = None
239      best_threshold = None
240
241      for feature_idx in range(X_train.shape[1]):
242          feature_values = X_train[:, feature_idx]
243          thresholds = np.unique(feature_values)
244
245          for threshold in thresholds:
246              predicted = (feature_values >= threshold).astype(int)
247              accuracy = np.mean(predicted == y_train)
```

```python
                if accuracy > best_accuracy:
                    best_accuracy = accuracy
                    best_feature = feature_idx
                    best_threshold = threshold

        best_feature_values = X_train[:, best_feature]
        predictions = (best_feature_values >= best_threshold).astype(int)

        return predictions, best_feature, best_threshold


y_train_pred, best_feature, best_threshold = oneR_classifier(X_train, y_train)
best_feature_values_val = X_val[:, best_feature]
y_val_pred = (best_feature_values_val >= best_threshold).astype(int)


accuracy_val = accuracy(y_val, y_val_pred)
conf_matrix_val = confusion_matrix_manual(y_val, y_val_pred)
precision_val = precision_manual(y_val, y_val_pred)
recall_val = recall_manual(y_val, y_val_pred)
f1_val = f1_score_manual(y_val, y_val_pred)

print("\nValidation Set Metrics:")
print(f"Accuracy: {accuracy_val}")
print(f"Confusion Matrix (TP, TN, FP, FN): {conf_matrix_val}")
print(f"Precision: {precision_val}")
print(f"Recall: {recall_val}")
print(f"F1-Score: {f1_val}")

y_val_prob = y_val_pred
precision_vals, recall_vals, _ = precision_recall_curve_manual(y_val, y_val_prob)
plot_precision_recall_curve(precision_vals, recall_vals, label="OneR Classifier (Validation)")

best_feature_values_test = X_test[:, best_feature]
y_test_pred = (best_feature_values_test >= best_threshold).astype(int)

accuracy_test = accuracy(y_test, y_test_pred)
conf_matrix_test = confusion_matrix_manual(y_test, y_test_pred)
precision_test = precision_manual(y_test, y_test_pred)
recall_test = recall_manual(y_test, y_test_pred)
f1_test = f1_score_manual(y_test, y_test_pred)

print("\nTest Set Metrics:")
print(f"Accuracy: {accuracy_test}")
print(f"Confusion Matrix (TP, TN, FP, FN): {conf_matrix_test}")
print(f"Precision: {precision_test}")
print(f"Recall: {recall_test}")
print(f"F1-Score: {f1_test}")


conf_matrix = confusion_matrix(y_test, y_test_pred)
print("Confusion Matrix:")
print(conf_matrix)

plt.figure(figsize=(8, 6))
sns.heatmap(conf_matrix, annot=True, fmt="d", cmap="Blues", xticklabels=[1, 2, 3], yticklabels
    =[1, 2, 3])
plt.title("Confusion Matrix")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()
```

```python
310 y_test_prob = y_test_pred
311 precision_vals, recall_test, _ = precision_recall_curve_manual(y_test, y_test_prob)
312 plot_precision_recall_curve(precision_vals, recall_test, label="OneR Classifier (Test)")
313
314 from sklearn.neighbors import KNeighborsClassifier
315 knn = KNeighborsClassifier(n_neighbors=3)
316
317
318 knn.fit(X_train, y_train)
319 y_val_pred = knn.predict(X_val)
320
321 accuracy_val = accuracy(y_val, y_val_pred)
322 precision_val = precision_manual(y_val, y_val_pred, average='weighted')
323 recall_val = recall_manual(y_val, y_val_pred, average='weighted')
324 f1_val = f1_score_manual(y_val, y_val_pred, average='weighted')
325
326
327 print(f"Validation Accuracy: {accuracy_val}")
328 print(f"Validation Precision: {precision_val}")
329 print(f"Validation Recall: {recall_val}")
330 print(f"Validation F1 Score: {f1_val}")
331 y_test_pred = knn.predict(X_test)
332 accuracy_test = accuracy(y_test, y_test_pred)
333 precision_test = precision_manual(y_test, y_test_pred, average='weighted')
334 recall_test = recall_manual(y_test, y_test_pred, average='weighted')
335 f1_test = f1_score_manual(y_test, y_test_pred, average='weighted')
336 print(f"Test Accuracy: {accuracy_test}")
337 print(f"Test Precision: {precision_test}")
338 print(f"Test Recall: {recall_test}")
339 print(f"Test F1 Score: {f1_test}")
340
341
342 conf_matrix = confusion_matrix(y_test, y_test_pred)
343 print("Confusion Matrix:")
344 print(conf_matrix)
345
346
347 plt.figure(figsize=(8, 6))
348 sns.heatmap(conf_matrix, annot=True, fmt="d", cmap="Blues", xticklabels=[1, 2, 3], yticklabels
        =[1, 2, 3])
349 plt.title("Confusion Matrix")
350 plt.xlabel("Predicted")
351 plt.ylabel("Actual")
352 plt.show()
353
354
355 from sklearn.metrics import roc_curve, auc
356
357
358 fpr, tpr, thresholds = roc_curve(y_test, knn.predict_proba(X_test)[:, 1], pos_label=1)
359 roc_auc = auc(fpr, tpr)
360
361 plt.figure(figsize=(8, 6))
362 plt.plot(fpr, tpr, color='blue', lw=2, label=f"ROC curve (area = {roc_auc:.2f})")
363 plt.plot([0, 1], [0, 1], color='gray', lw=2, linestyle='--')
364 plt.xlabel('False Positive Rate')
365 plt.ylabel('True Positive Rate')
366 plt.title('Receiver Operating Characteristic (ROC) Curve')
367 plt.legend(loc='lower right')
368 plt.show()
369
370
371 param_grid = {
```

```python
372        'n_neighbors': [3, 5, 7, 9, 11],
373        'weights': ['uniform', 'distance'],
374        'metric': ['euclidean', 'manhattan']
375    }
376
377    grid_search = GridSearchCV(KNeighborsClassifier(), param_grid, cv=5, scoring='accuracy')
378    grid_search.fit(X_train, y_train)
379
380    print("Best Parameters from Grid Search:")
381    print(grid_search.best_params_)
382    print(f"Best Accuracy from Grid Search: {grid_search.best_score_}")
383
384    best_knn = grid_search.best_estimator_
385
386    y_test_pred_best = best_knn.predict(X_test)
387
388    accuracy_best = accuracy(y_test, y_test_pred_best)
389    precision_best = precision_manual(y_test, y_test_pred_best, average='weighted')
390    recall_best = recall_manual(y_test, y_test_pred_best, average='weighted')
391    f1_best = f1_score_manual(y_test, y_test_pred_best, average='weighted')
392
393    print(f"Test Accuracy (Best K-NN): {accuracy_best}")
394    print(f"Test Precision (Best K-NN): {precision_best}")
395    print(f"Test Recall (Best K-NN): {recall_best}")
396    print(f"Test F1 Score (Best K-NN): {f1_best}")
397
398
399    k_values = [3, 5, 7, 9, 11]
400    accuracies = []
401
402    for k in k_values:
403        knn = KNeighborsClassifier(n_neighbors=k)
404        knn.fit(X_train, y_train)
405        y_pred = knn.predict(X_test)
406        accuracies.append(accuracy(y_test, y_pred))
407
408    plt.figure(figsize=(8, 6))
409    plt.plot(k_values, accuracies, marker='o', color='green')
410    plt.title('Accuracy vs. Number of Neighbors (K)')
411    plt.xlabel('Number of Neighbors (K)')
412    plt.ylabel('Accuracy')
413    plt.grid(True)
414    plt.show()
415
416    k_values = [3, 5, 7, 9, 11]
417    accuracies = []
418
419    for k in k_values:
420        knn = KNeighborsClassifier(n_neighbors=k)
421        knn.fit(X_train, y_train)
422        y_val_pred = knn.predict(X_val)
423        accuracies.append(accuracy(y_val, y_val_pred))
424
425    plt.figure(figsize=(8, 6))
426    plt.plot(k_values, accuracies, marker='o', color='green')
427    plt.title('Accuracy vs. Number of Neighbors (K)')
428    plt.xlabel('Number of Neighbors (K)')
429    plt.ylabel('Accuracy')
430    plt.grid(True)
431    plt.show()
432
433
434
```

```python
435  from sklearn.naive_bayes import GaussianNB
436
437
438  nb = GaussianNB()
439  nb.fit(X_train, y_train)
440  y_val_pred = nb.predict(X_val)
441  y_test_pred = nb.predict(X_test)
442  y_test_pred = nb.predict(X_test)
443
444
445  accuracy_val = accuracy(y_val, y_val_pred)
446  precision_val = precision_manual(y_val, y_val_pred, average='weighted')
447  recall_val = recall_manual(y_val, y_val_pred, average='weighted')
448  f1_val = f1_score_manual(y_val, y_val_pred, average='weighted')
449
450
451
452  print(f"Validation Accuracy: {accuracy_val}")
453  print(f"Validation Precision: {precision_val}")
454  print(f"Validation Recall: {recall_val}")
455  print(f"Validation F1 Score: {f1_val}")
456
457
458  accuracy_test = accuracy(y_test, y_test_pred)
459  precision_test = precision_manual(y_test, y_test_pred, average='weighted')
460  recall_test = recall_manual(y_test, y_test_pred, average='weighted')
461  f1_test = f1_score_manual(y_test, y_test_pred, average='weighted')
462
463
464  print(f"Test Accuracy: {accuracy_test}")
465  print(f"Test Precision: {precision_test}")
466  print(f"Test Recall: {recall_test}")
467  print(f"Test F1 Score: {f1_test}")
468
469  conf_matrix = confusion_matrix(y_test, y_test_pred)
470  print("Confusion Matrix:")
471  print(conf_matrix)
472
473  plt.figure(figsize=(8, 6))
474  sns.heatmap(conf_matrix, annot=True, fmt="d", cmap="Blues", xticklabels=[1, 2, 3], yticklabels
          =[1, 2, 3])
475  plt.title("Confusion Matrix")
476  plt.xlabel("Predicted")
477  plt.ylabel("Actual")
478  plt.show()
479
480
481
482  accuracy_test = accuracy(y_test, y_test_pred)
483  precision_test = precision_manual(y_test, y_test_pred, average='weighted')
484  recall_test = recall_manual(y_test, y_test_pred, average='weighted')
485  f1_test = f1_score_manual(y_test, y_test_pred, average='weighted')
486
487  print(f"Test Accuracy (Naive Bayes): {accuracy_test}")
488  print(f"Test Precision (Naive Bayes): {precision_test}")
489  print(f"Test Recall (Naive Bayes): {recall_test}")
490  print(f"Test F1 Score (Naive Bayes): {f1_test}")
491
492
493
494  svm.fit(X_train, y_train)
495  y_test_pred_svm = svm.predict(X_test)
496
```

```
497
498
499
500 accuracy_test_svm = accuracy(y_test, y_test_pred_svm)
501 precision_test_svm = precision_manual(y_test, y_test_pred_svm, average='weighted')
502 recall_test_svm = recall_manual(y_test, y_test_pred_svm, average='weighted')
503 f1_test_svm = f1_score_manual(y_test, y_test_pred_svm, average='weighted')
504
505 print(f"Test Accuracy (SVM): {accuracy_test_svm}")
506 print(f"Test Precision (SVM): {precision_test_svm}")
507 print(f"Test Recall (SVM): {recall_test_svm}")
508 print(f"Test F1 Score (SVM): {f1_test_svm}")
509
510
511
512 conf_matrix_svm = confusion_matrix(y_test, y_test_pred_svm)
513
514 plt.figure(figsize=(8, 6))
515 sns.heatmap(conf_matrix_svm, annot=True, fmt="d", cmap="Blues", xticklabels=[1, 2, 3],
        yticklabels=[1, 2, 3])
516 plt.title("SVM Confusion Matrix")
517 plt.xlabel("Predicted")
518 plt.ylabel("Actual")
519 plt.show()
520
521
522 y_test_bin = label_binarize(y_test, classes=[1, 2, 3])
523 y_score_svm = svm.decision_function(X_test)
524
525 print(f"Shape of y_score_svm: {y_score_svm.shape}")
526 print(f"Shape of y_test_bin: {y_test_bin.shape}")
527
528
529
530 y_score_svm = svm.predict_proba(X_test)[:, 1]
531
532 precision, recall, _ = precision_recall_curve(y_test, y_score_svm)
533
534 plt.figure(figsize=(8, 6))
535 plt.plot(recall, precision, marker='.', label='SVM')
536 plt.title('Precision-Recall Curve for SVM (Binary Classification)')
537 plt.xlabel('Recall')
538 plt.ylabel('Precision')
539 plt.legend(loc='best')
540 plt.show()
541
542
543 y_score_svm = svm.predict_proba(X_test)[:, 1]
544
545 precision, recall, _ = precision_recall_curve(y_test, y_score_svm)
546
547 plt.figure(figsize=(8, 6))
548 plt.plot(recall, precision, marker='.', label='SVM')
549 plt.title('Precision-Recall Curve for SVM (Binary Classification)')
550 plt.xlabel('Recall')
551 plt.ylabel('Precision')
552 plt.legend(loc='best')
553 plt.show()
```

## B. Car Dataset

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

```python
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.metrics import precision_recall_curve, average_precision_score
from sklearn.preprocessing import label_binarize
import sklearn
import warnings
%matplotlib inline
sns.set()
warnings.filterwarnings('ignore')
df = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/car.data")
df.head()
df.columns = ["buying", "maint", "doors", "persons", "lug_boot", "safety", "class values"]

df['doors'] = df['doors'].replace(['5more'], '5')
df['persons'] = df['persons'].replace(['more'], '4')
for col in ['buying', 'maint', 'doors', 'persons', 'lug_boot', 'safety']:
    df[col] = df[col].astype('category')

for col in ['buying', 'maint', 'doors', 'persons', 'lug_boot', 'safety']:
    df[col] = df[col].cat.codes
df = pd.get_dummies(df, columns=['buying', 'maint', 'doors', 'persons', 'lug_boot', 'safety'])

df.head()

df['class values'].unique()
df = df.drop(['buying_3', 'maint_3', 'doors_3', 'persons_1', 'lug_boot_2', 'safety_2'], axis=1)
df.head()
from sklearn.preprocessing import LabelEncoder
labelencoder = LabelEncoder()
df['class values'] = labelencoder.fit_transform(df['class values'])
x = df.drop(["class values"], axis=1)
y = df["class values"]
x.head()
y.unique()
x_train, x_temp, y_train, y_temp = train_test_split(x, y, test_size=0.3, random_state=42)
x_val, x_test, y_val, y_test = train_test_split(x_temp, y_temp, test_size=0.5, random_state=42)

class OneRClassifier:
    def __init__(self):
        self.rule = {}
        self.best_feature = None

    def fit(self, X, y):
        best_error = float('inf')

        for feature in X.columns:
            feature_rule = {}
            grouped = pd.DataFrame({'feature': X[feature], 'target': y}).groupby('feature')['target'].agg(lambda x: x.value_counts().idxmax())
            feature_rule = grouped.to_dict()

            y_pred = X[feature].map(feature_rule)
            error = sum(y_pred != y)

            if error < best_error:
                best_error = error
                self.rule = feature_rule
                self.best_feature = feature

    def predict(self, X):
```

```python
            if not self.best_feature:
                raise ValueError("The model has not been trained yet!")
            return X[self.best_feature].map(self.rule)

    def calculate_confusion_matrix(y_true, y_pred, num_classes):
        cm = np.zeros((num_classes, num_classes), dtype=int)
        for true, pred in zip(y_true, y_pred):
            cm[true][pred] += 1
        return cm

    def calculate_metrics(y_true, y_pred):
        num_classes = len(np.unique(y_true))
        cm = calculate_confusion_matrix(y_true, y_pred, num_classes)
        total = cm.sum()
        accuracy = np.trace(cm) / total

        precisions = []
        recalls = []
        f1_scores = []
        for i in range(num_classes):
            tp = cm[i, i]
            fp = sum(cm[:, i]) - tp
            fn = sum(cm[i, :]) - tp
            precision = tp / (tp + fp) if (tp + fp) > 0 else 0
            recall = tp / (tp + fn) if (tp + fn) > 0 else 0
            f1 = 2 * (precision * recall) / (precision + recall) if (precision + recall) > 0 else 0

            precisions.append(precision)
            recalls.append(recall)
            f1_scores.append(f1)

        avg_precision = np.mean(precisions)
        avg_recall = np.mean(recalls)
        avg_f1 = np.mean(f1_scores)

        return {
            'accuracy': accuracy,
            'confusion_matrix': cm,
            'average_precision': avg_precision,
            'average_recall': avg_recall,
            'average_f1': avg_f1,
        }

model = OneRClassifier()
model.fit(x_train, y_train)

y_train_pred = model.predict(x_train)
y_val_pred = model.predict(x_val)
y_test_pred = model.predict(x_test)

metrics = calculate_metrics(y_val.to_numpy(), y_val_pred.to_numpy())
metrics2 = calculate_metrics(y_test.to_numpy(), y_test_pred.to_numpy())

print("OneR Classifier Validation Metrics:")
print(f"Accuracy: {metrics['accuracy']:.2f}")
print("Confusion Matrix (Custom):")
print(metrics['confusion_matrix'])
print(f"Average Precision: {metrics['average_precision']:.2f}")
print(f"Average Recall: {metrics['average_recall']:.2f}")
print(f"Average F1-Score: {metrics['average_f1']:.2f}")
def plot_confusion_matrix(cm, class_names):
    plt.figure(figsize=(8, 6))
```

```python
127         sns.heatmap(cm, annot=True, fmt='d', cmap="Blues", xticklabels=class_names, yticklabels=
            class_names)
128         plt.title("Confusion Matrix")
129         plt.xlabel("Predicted Class")
130         plt.ylabel("True Class")
131         plt.show()
132
133 class_names = labelencoder.inverse_transform(range(len(np.unique(y_train))))
134
135 plot_confusion_matrix(metrics['confusion_matrix'], class_names)
136 print("OneR Classifier test Metrics:")
137 print(f"Accuracy: {metrics2['accuracy']:.2f}")
138 print("Confusion Matrix (Custom):")
139 print(metrics2['confusion_matrix'])
140 print(f"Average Precision: {metrics2['average_precision']:.2f}")
141 print(f"Average Recall: {metrics2['average_recall']:.2f}")
142 print(f"Average F1-Score: {metrics2['average_f1']:.2f}")
143 def plot_confusion_matrix(cm, class_names):
144     plt.figure(figsize=(8, 6))
145     sns.heatmap(cm, annot=True, fmt='d', cmap="Blues", xticklabels=class_names, yticklabels=
            class_names)
146     plt.title("Confusion Matrix")
147     plt.xlabel("Predicted Class")
148     plt.ylabel("True Class")
149     plt.show()
150
151 class_names = labelencoder.inverse_transform(range(len(np.unique(y_train))))
152
153 plot_confusion_matrix(metrics2['confusion_matrix'], class_names)
154 majority_class = y_train.mode()[0]
155 y_val_pred_zeror = [majority_class] * len(y_val)
156 y_test_pred_zeror = [majority_class] * len(y_test)
157 def accuracy_score(y_true, y_pred):
158     y_true = np.array(y_true)
159     y_pred = np.array(y_pred)
160
161     corr_pred = np.sum(y_true == y_pred)
162
163     total_pred = len(y_true)
164
165     accuracy = corr_pred / total_pred
166
167     return accuracy
168 def confusion_matrix(y_true, y_pred, labels=None):
169     y_true = np.array(y_true)
170     y_pred = np.array(y_pred)
171
172     if labels is None:
173         labels = np.unique(np.concatenate([y_true, y_pred]))
174
175     cm = np.zeros((len(labels), len(labels)), dtype=int)
176
177     label_to_index = {label: idx for idx, label in enumerate(labels)}
178
179     for true_label, pred_label in zip(y_true, y_pred):
180         true_idx = label_to_index[true_label]
181         pred_idx = label_to_index[pred_label]
182         cm[true_idx, pred_idx] += 1
183
184     return cm
185
186 def plot_confusion_matrix(cm, labels):
187     plt.figure(figsize=(8, 6))
```

```python
188        plt.imshow(cm, interpolation='nearest', cmap=plt.cm.Blues)
189        plt.title('Confusion Matrix')
190        plt.colorbar()
191
192        tick_marks = np.arange(len(labels))
193        plt.xticks(tick_marks, labels, rotation=45)
194        plt.yticks(tick_marks, labels)
195        for i in range(len(labels)):
196            for j in range(len(labels)):
197                plt.text(j, i, cm[i, j],
198                         horizontalalignment="center",
199                         color="white" if cm[i, j] > cm.max() / 2 else "black")
200
201        plt.tight_layout()
202        plt.ylabel('True label')
203        plt.xlabel('Predicted label')
204        plt.show()
205    def precision_recall_f1(cm):
206        epsilon = 1e-10
207        precision = np.diag(cm) / (np.sum(cm, axis=0) + epsilon)
208        recall = np.diag(cm) / (np.sum(cm, axis=1) + epsilon)
209        f1 = 2 * (precision * recall) / (precision + recall + epsilon)
210
211        accuracy = np.sum(np.diag(cm)) / np.sum(cm)
212        f1 = np.nan_to_num(f1)
213
214        return precision, recall, f1, accuracy
215
216
217    def classification_report(y_true, y_pred, labels=None):
218        if labels is None:
219            labels = np.unique(np.concatenate([y_true, y_pred]))
220
221        cm = confusion_matrix(y_true, y_pred, labels)
222
223        precision, recall, f1, accuracy = precision_recall_f1(cm)
224
225        report = "                    precision    recall  f1-score   support\n\n"
226
227        for idx, label in enumerate(labels):
228            report += f"{label:>10}        {precision[idx]:.2f}        {recall[idx]:.2f}        {f1[idx
       ]:.2f}        {np.sum(cm[idx]):>3}\n"
229
230        macro_avg_precision = np.mean(precision)
231        macro_avg_recall = np.mean(recall)
232        macro_avg_f1 = np.mean(f1)
233        report += f"\n    accuracy                                {accuracy:.2f}        {len(y_true)}\n"
234        report += f"    macro avg        {macro_avg_precision:.2f}        {macro_avg_recall:.2f}        {
       macro_avg_f1:.2f}        {len(y_true)}\n"
235
236        return report
237    val_accuracy = accuracy_score(y_val, y_val_pred_zeror)
238    test_accuracy = accuracy_score(y_test, y_test_pred_zeror)
239
240    print(f"ZeroR Validation Accuracy (Custom): {val_accuracy:.2f}")
241    print(f"ZeroR Test Accuracy (Custom): {test_accuracy:.2f}")
242    labels = np.unique(y_test)
243    cm = confusion_matrix(y_test, y_test_pred_zeror, labels=labels)
244
245    report = classification_report(y_test, y_test_pred_zeror, labels=labels)
246    print("\nCustom Classification Report (Test Data):")
247    print(report)
248    labels = np.unique(y_val)
```

```python
249  cm = confusion_matrix(y_val, y_val_pred_zeror, labels=labels)
250
251  print("Confusion Matrix:")
252  print(cm)
253
254  plot_confusion_matrix(cm, labels=labels)
255  labels = np.unique(y_test)
256  cm = confusion_matrix(y_test, y_test_pred_zeror, labels=labels)
257
258  print("Confusion Matrix:")
259  print(cm)
260
261  plot_confusion_matrix(cm, labels=labels)
262
263  from sklearn.neighbors import KNeighborsClassifier
264  from sklearn.metrics import  PrecisionRecallDisplay
265  from sklearn.model_selection import GridSearchCV
266
267  param_grid = {'n_neighbors': range(1, 21), 'weights': ['uniform', 'distance'], 'metric': ['
           euclidean', 'manhattan']}
268  grid_search = GridSearchCV(KNeighborsClassifier(), param_grid, cv=5, scoring='accuracy')
269  grid_search.fit(x_train, y_train)
270
271  best_knn = grid_search.best_estimator_
272  print("Best Parameters:", grid_search.best_params_)
273
274  y_test_pred = best_knn.predict(x_test)
275  y_val_pred = best_knn.predict(x_val)
276
277  val_accuracy = accuracy_score(y_val, y_val_pred)
278  print("Validation Accuracy:", val_accuracy)
279
280  test_accuracy = accuracy_score(y_test, y_test_pred)
281  print("Test Accuracy:", test_accuracy)
282  cm1 = confusion_matrix(y_val, y_val_pred)
283  print("Validation Confusion Matrix:")
284  print(cm1)
285
286  cm2 = confusion_matrix(y_test, y_test_pred)
287  print("Test Confusion Matrix:")
288  print(cm2)
289
290  def plot_confusion_matrix(cm, labels):
291      plt.figure(figsize=(8, 6))
292      plt.imshow(cm, interpolation='nearest', cmap=plt.cm.Blues)
293      plt.title('Confusion Matrix')
294      plt.colorbar()
295
296      tick_marks = np.arange(len(labels))
297      plt.xticks(tick_marks, labels, rotation=45)
298      plt.yticks(tick_marks, labels)
299
300      for i in range(len(labels)):
301          for j in range(len(labels)):
302              plt.text(j, i, cm[i, j],
303                       horizontalalignment="center",
304                       color="white" if cm[i, j] > cm.max() / 2 else "black")
305
306      plt.tight_layout()
307      plt.ylabel('True label')
308      plt.xlabel('Predicted label')
309      plt.show()
310  plot_confusion_matrix(cm1, labels=labelencoder.classes_)
```

```
311
312 plot_confusion_matrix(cm2, labels=labelencoder.classes_)
313
314 val_report = classification_report(y_val, y_val_pred)
315 print("Validation Classification Report:")
316 print(val_report)
317
318 report = classification_report(y_test, y_test_pred)
319 print("Test Classification Report:")
320 print(report)
321 y_val_pred_prob = knn.predict_proba(x_val)
322
323 for i in range(len(labelencoder.classes_)):
324     precision, recall, _ = precision_recall_curve(y_val == i, y_val_pred_prob[:, i])
325     PrecisionRecallDisplay(precision=precision, recall=recall, average_precision=np.mean(
    precision)).plot()
326     plt.title(f'Validation Precision-Recall Curve for class {labelencoder.classes_[i]}')
327
328 plt.show()
329 y_test_pred_prob = knn.predict_proba(x_test)
330
331 for i in range(len(labelencoder.classes_)):
332     precision, recall, _ = precision_recall_curve(y_test == i, y_test_pred_prob[:, i])
333     PrecisionRecallDisplay(precision=precision, recall=recall, average_precision=np.mean(
    precision)).plot()
334     plt.title(f'Test Precision-Recall Curve for class {labelencoder.classes_[i]}')
335
336 plt.show()
337 from sklearn.naive_bayes import GaussianNB
338 from sklearn.metrics import PrecisionRecallDisplay
339 y_test_pred = nb.predict(x_test)
340 y_val_pred = nb.predict(x_val)
341 accuracy = accuracy_score(y_test, y_test_pred)
342 print("Test Accuracy:", accuracy)
343 accuracy = accuracy_score(y_val, y_val_pred)
344 print("Validation Accuracy:", accuracy)
345 cm1 = confusion_matrix(y_test, y_test_pred)
346 print("Test Confusion Matrix:")
347 print(cm1)
348
349 cm2 = confusion_matrix(y_val, y_val_pred)
350 print("Validation Confusion Matrix:")
351 print(cm2)
352 def plot_confusion_matrix(cm, labels):
353     plt.figure(figsize=(8, 6))
354     plt.imshow(cm, interpolation='nearest', cmap=plt.cm.Blues)
355     plt.title('Confusion Matrix')
356     plt.colorbar()
357
358     tick_marks = np.arange(len(labels))
359     plt.xticks(tick_marks, labels, rotation=45)
360     plt.yticks(tick_marks, labels)
361
362     for i in range(len(labels)):
363         for j in range(len(labels)):
364             plt.text(j, i, cm[i, j],
365                     horizontalalignment="center",
366                     color="white" if cm[i, j] > cm.max() / 2 else "black")
367
368     plt.tight_layout()
369     plt.ylabel('True label')
370     plt.xlabel('Predicted label')
371     plt.show()
```

```python
372  plot_confusion_matrix(cm1, labels=labelencoder.classes_)
373  plot_confusion_matrix(cm2, labels=labelencoder.classes_)
374  test_report = classification_report(y_test, y_test_pred)
375  print("Test Classification Report:")
376  print(test_report)
377  val_report = classification_report(y_val, y_val_pred)
378  print("Validation Classification Report:")
379  print(val_report)
380  y_test_pred_prob = nb.predict_proba(x_test)
381
382  for i in range(len(labelencoder.classes_)):
383      precision, recall, _ = precision_recall_curve(y_test == i, y_test_pred_prob[:, i])
384      PrecisionRecallDisplay(precision=precision, recall=recall, average_precision=np.mean(
         precision)).plot()
385      plt.title(f'Precision-Recall Curve for class {labelencoder.classes_[i]}')
386
387  plt.show()
388
389  y_val_pred_prob = nb.predict_proba(x_val)
390
391  for i in range(len(labelencoder.classes_)):
392      precision, recall, _ = precision_recall_curve(y_val == i, y_val_pred_prob[:, i])
393      PrecisionRecallDisplay(precision=precision, recall=recall, average_precision=np.mean(
         precision)).plot()
394      plt.title(f'Precision-Recall Curve for class {labelencoder.classes_[i]}')
395
396  plt.show()
397
398  from sklearn.svm import SVC
399  from sklearn.metrics import PrecisionRecallDisplay
400
401
402  # svm = SVC(kernel='linear', probability=True)
403  # svm.fit(x_train, y_train)
404  from sklearn.model_selection import GridSearchCV
405  param_grid = {
406      'C': [0.1, 1, 10],
407      'kernel': ['linear', 'rbf'],
408      'gamma': [0.001, 0.01, 0.1]
409  }
410  grid_search = GridSearchCV(SVC(probability=True), param_grid, cv=5)
411  grid_search.fit(x_train, y_train)
412  best_model = grid_search.best_estimator_
413  # y_val_pred = svm.predict(x_val)
414  y_val_pred = best_model.predict(x_val)
415
416  val_accuracy = accuracy_score(y_val, y_val_pred)
417  print("Validation Accuracy:", val_accuracy)
418
419  val_cm = confusion_matrix(y_val, y_val_pred)
420  print("Validation Confusion Matrix:")
421  print(val_cm)
422
423  plot_confusion_matrix(val_cm, labels=labelencoder.classes_)
424
425  val_report = classification_report(y_val, y_val_pred)
426  print("Validation Classification Report:")
427  print(val_report)
428  y_val_pred_prob = best_model.predict_proba(x_val)
429  for i in range(len(labelencoder.classes_)):
430      precision, recall, _ = precision_recall_curve(y_val == i, y_val_pred_prob[:, i])
431      PrecisionRecallDisplay(precision=precision, recall=recall, average_precision=np.mean(
         precision)).plot()
```

```
432        plt.title(f'Precision-Recall Curve for class {labelencoder.classes_[i]} (Validation)')

434  plt.show()

436  # y_test_pred = svm.predict(x_test)
437  y_test_pred = best_model.predict(x_test)
438  test_accuracy = accuracy_score(y_test, y_test_pred)
439  print("Test Accuracy:", test_accuracy)
440  test_cm = confusion_matrix(y_test, y_test_pred)
441  print("Test Confusion Matrix:")
442  print(test_cm)

444  plot_confusion_matrix(test_cm, labels=labelencoder.classes_)

446  test_report = classification_report(y_test, y_test_pred)
447  print("Test Classification Report:")
448  print(test_report)

450  y_test_pred_prob = best_model.predict_proba(x_test)
451  for i in range(len(labelencoder.classes_)):
452      precision, recall, _ = precision_recall_curve(y_test == i, y_test_pred_prob[:, i])
453      PrecisionRecallDisplay(precision=precision, recall=recall, average_precision=np.mean(
         precision)).plot()
454      plt.title(f'Precision-Recall Curve for class {labelencoder.classes_[i]} (Test)')

456  plt.show()
```

## C. 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  import matplotlib.pyplot as plt
11  df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/DT-Wage.csv')
12  from sklearn.preprocessing import StandardScaler

14  df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/DT-Wage.csv')

16  df.shape

18  df.info()
19  un = df['maritl'].unique()
20  print(un)

22  un = df['race'].unique()
23  print(un)

25  un = df['education'].unique()
26  print(un)

28  un = df['region'].unique()
29  print(un)

31  un = df['jobclass'].unique()
32  print(un)

34  un = df['health'].unique()
```

```python
35 print(un)
36
37 categorical = df.select_dtypes(include=['object']).columns
38 numerical = df.select_dtypes(exclude=['object']).columns
39
40 encoder = OneHotEncoder(drop='first', sparse_output=False)
41 encoded_data = encoder.fit_transform(df[categorical])
42 encoded_df = pd.DataFrame(encoded_data, columns=encoder.get_feature_names_out(categorical))
43 encoded_df.index = df.index
44
45 df = df.drop(categorical, axis=1)
46 df = pd.concat([df, encoded_df], axis=1)
47
48 X = df.drop(columns=['wage'])
49 y = df['wage']
50 X.head()
51
52 X.shape
53
54 scaler = StandardScaler()
55 X_scaled = scaler.fit_transform(X)
56
57 X_train, X_temp, y_train, y_temp = train_test_split(X_scaled, y, test_size=0.3, random_state
       =42)
58 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)
59
60 from sklearn.svm import SVR
61 from sklearn.model_selection import GridSearchCV
62 from sklearn.metrics import mean_squared_error, r2_score
63 svr = SVR()
64
65 param_grid = {
66     'C': [0.1, 1, 10, 100],
67     'kernel': ['linear', 'poly', 'rbf'],
68     'degree': [3, 4, 5],
69     'gamma': ['scale', 'auto'],
70     'epsilon': [0.1, 0.2, 0.5],
71 }
72
73 grid_search = GridSearchCV(svr, param_grid, cv=5, scoring='neg_mean_squared_error', n_jobs=-1)
74 grid_search.fit(X_train, y_train)
75 print("Best hyperparameters:", grid_search.best_params_)
76
77 best_svr = grid_search.best_estimator_
78 y_pred_val = best_svr.predict(X_val)
79 def custom_mse(y_true, y_pred):
80     return np.mean((y_true - y_pred) ** 2)
81
82 def custom_r2(y_true, y_pred):
83     ss_total = np.sum((y_true - np.mean(y_true)) ** 2)
84     ss_residual = np.sum((y_true - y_pred) ** 2)
85     return 1 - (ss_residual / ss_total)
86
87 mse = custom_mse(y_val, y_pred_val)
88 r2 = custom_r2(y_val, y_pred_val)
89 print(f"Validation MSE: {mse}")
90 print(f"Validation R^2: {r2}")
91
92
93 y_pred_test = best_svr.predict(X_test)
94 test_mse = custom_mse(y_test, y_pred_test)
95 test_r2 = custom_r2(y_test, y_pred_test)
96 print(f"Test MSE: {test_mse}")
```

```python
97  print(f"Test R^2: {test_r2}")
98
99  plt.figure(figsize=(10, 6))
100 plt.scatter(y_val, y_pred_val, color='blue', alpha=0.5)
101 plt.plot([min(y_val), max(y_val)], [min(y_val), max(y_val)], color='red', linestyle='--')
102 plt.title("SVR: Actual vs Predicted (Validation Set)")
103 plt.xlabel("Actual Values")
104 plt.ylabel("Predicted Values")
105 plt.show()
106
107 residuals = y_val - y_pred_val
108 plt.figure(figsize=(10, 6))
109 plt.scatter(y_pred_val, residuals, color='green', alpha=0.5)
110 plt.hlines(0, xmin=min(y_pred_val), xmax=max(y_pred_val), colors='red', linestyle='--')
111 plt.title("Residuals Plot (Validation Set)")
112 plt.xlabel("Predicted Values")
113 plt.ylabel("Residuals")
114 plt.show()
115
116 plt.figure(figsize=(10, 6))
117 plt.scatter(y_test, y_pred_test, color='blue', alpha=0.5)
118 plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)], color='red', linestyle='--')
119 plt.title("SVR: Actual vs Predicted (Test Set)")
120 plt.xlabel("Actual Values")
121 plt.ylabel("Predicted Values")
122 plt.show()
123
124 residuals_test = y_test - y_pred_test
125 plt.figure(figsize=(10, 6))
126 plt.scatter(y_pred_test, residuals_test, color='green', alpha=0.5)
127 plt.hlines(0, xmin=min(y_pred_test), xmax=max(y_pred_test), colors='red', linestyle='--')
128 plt.title("Residuals Plot (Test Set)")
129 plt.xlabel("Predicted Values")
130 plt.ylabel("Residuals")
131 plt.show()
```

### D. Credit Dataset

```python
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 import matplotlib.pyplot as plt
11 from sklearn.preprocessing import StandardScaler
12 df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/DT-Credit.csv')
13 df.head()
14 df.shape
15 df.describe()
16 df.isnull().sum()
17 un = df['Own'].unique()
18 print(un)
19
20 un = df['Student'].unique()
21 print(un)
22
23 un = df['Region'].unique()
24 print(un)
25
```

```python
26 un = df['Married'].unique()
27 print(un)
28
29 categorical = df.select_dtypes(include=['object']).columns
30 numerical = df.select_dtypes(exclude=['object']).columns
31
32 encoder = OneHotEncoder(drop='first', sparse_output=False)
33 encoded_data = encoder.fit_transform(df[categorical])
34 encoded_df = pd.DataFrame(encoded_data, columns=encoder.get_feature_names_out(categorical))
35 encoded_df.index = df.index
36
37 df = df.drop(categorical, axis=1)
38 df = pd.concat([df, encoded_df], axis=1)
39 X = df.drop(columns=['Balance'])
40 y = df['Balance']
41 X.head()
42
43 X.shape
44
45 corr = df.corr()
46 target = corr['Balance'].sort_values(ascending=False)
47 print("Correlation of each feature with the target variable (Balance):")
48 print(target)
49
50 fig, ax = plt.subplots(1,1, figsize=(12, 12))
51 df.boxplot('Income', 'Balance', ax=ax)
52 plt.suptitle('Income vs Balance')
53 plt.title('')
54 plt.ylabel('Income')
55 plt.xticks(rotation=90)
56 plt.show()
57
58 fig, ax = plt.subplots(1,1, figsize=(12, 12))
59 df.boxplot('Age', 'Balance', ax=ax)
60 plt.suptitle('Age vs Balance')
61 plt.title('')
62 plt.ylabel('Age')
63 plt.xticks(rotation=90)
64 plt.show()
65
66 scaler = StandardScaler()
67 X_scaled = scaler.fit_transform(X)
68 X_train, X_temp, y_train, y_temp = train_test_split(X_scaled, y, test_size=0.3, random_state
     =42)
69 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)
70
71 X_train.shape
72 X_val.shape
73 X_test.shape
74
75 from sklearn.model_selection import GridSearchCV
76 from sklearn.svm import SVR
77 from sklearn.metrics import mean_squared_error, r2_score
78
79 param_grid = {
80     'C': [0.1, 1, 10, 100,200,300,500,1500,2000],
81     'epsilon': [0.01, 0.1, 0.2,0.25,0.3],
82     'gamma': ['scale', 'auto']
83 }
84
85 svr = SVR(kernel='rbf')
86
87 grid_search = GridSearchCV(estimator=svr, param_grid=param_grid,
```

```
88                                cv=5, scoring='neg_mean_squared_error', n_jobs=-1)
89
90 grid_search.fit(X_train, y_train)
91 best_params = grid_search.best_params_
92 print("Best hyperparameters found by GridSearchCV:")
93 print(best_params)
94
95 best_svr = grid_search.best_estimator_
96 y_pred_train = best_svr.predict(X_train)
97 y_pred_val = best_svr.predict(X_val)
98 y_pred_test = best_svr.predict(X_test)
99 def custom_mse(y_true, y_pred):
100     return np.mean((y_true - y_pred) ** 2)
101
102 def custom_r2(y_true, y_pred):
103     ss_total = np.sum((y_true - np.mean(y_true)) ** 2)
104     ss_residual = np.sum((y_true - y_pred) ** 2)
105     return 1 - (ss_residual / ss_total)
106
107 mse_train = custom_mse(y_train, y_pred_train)
108 mse_val = custom_mse(y_val, y_pred_val)
109 mse_test = custom_mse(y_test, y_pred_test)
110
111 r2_train = custom_r2(y_train, y_pred_train)
112 r2_val = custom_r2(y_val, y_pred_val)
113 r2_test = custom_r2(y_test, y_pred_test)
114
115 print("SVR Model with GridSearchCV Evaluation:")
116 print(f"Training MSE: {mse_train}")
117 print(f"Validation MSE: {mse_val}")
118 print(f"Test MSE: {mse_test}")
119 print(f"Training R : {r2_train}")
120 print(f"Validation R : {r2_val}")
121 print(f"Test  R : {r2_test}")
122
123
124 plt.figure(figsize=(12, 4))
125 plt.subplot(1, 3, 1)
126 plt.scatter(y_train, y_pred_train, color='blue', alpha=0.5)
127 plt.plot([y_train.min(), y_train.max()], [y_train.min(), y_train.max()], color='red', lw=2)
128 plt.title('Training Set: Predicted vs Actual')
129 plt.xlabel('Actual')
130 plt.ylabel('Predicted')
131 plt.tight_layout()
132 plt.show()
133
134 plt.figure(figsize=(12, 4))
135 plt.subplot(1, 3, 2)
136 plt.scatter(y_val, y_pred_val, color='green', alpha=0.5)
137 plt.plot([y_val.min(), y_val.max()], [y_val.min(), y_val.max()], color='red', lw=2)
138 plt.title('Validation Set: Predicted vs Actual')
139 plt.xlabel('Actual')
140 plt.ylabel('Predicted')
141 plt.tight_layout()
142 plt.show()
143
144 plt.figure(figsize=(12, 4))
145 plt.subplot(1, 3, 3)
146 plt.scatter(y_test, y_pred_test, color='orange', alpha=0.5)
147 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red', lw=2)
148 plt.title('Test Set: Predicted vs Actual')
149 plt.xlabel('Actual')
150 plt.ylabel('Predicted')
```

```python
plt.tight_layout()
plt.show()
plt.figure(figsize=(12, 4))
plt.subplot(1, 3, 1)
plt.scatter(y_pred_train, y_pred_train - y_train, color='blue', alpha=0.5)
plt.hlines(y=0, xmin=y_pred_train.min(), xmax=y_pred_train.max(), color='red', lw=2)
plt.title('Training Set: Residuals')
plt.xlabel('Predicted')
plt.ylabel('Residuals (Error)')

plt.tight_layout()
plt.show()
plt.figure(figsize=(12, 4))
plt.subplot(1, 3, 2)
plt.scatter(y_pred_val, y_pred_val - y_val, color='green', alpha=0.5)
plt.hlines(y=0, xmin=y_pred_val.min(), xmax=y_pred_val.max(), color='red', lw=2)
plt.title('Validation Set: Residuals')
plt.xlabel('Predicted')
plt.ylabel('Residuals (Error)')

plt.tight_layout()
plt.show()
plt.figure(figsize=(12, 4))
plt.subplot(1, 3, 3)
plt.scatter(y_pred_test, y_pred_test - y_test, color='orange', alpha=0.5)
plt.hlines(y=0, xmin=y_pred_test.min(), xmax=y_pred_test.max(), color='red', lw=2)
plt.title('Test Set: Residuals')
plt.xlabel('Predicted')
plt.ylabel('Residuals (Error)')

plt.tight_layout()
plt.show()
```