

# Comprehensive Car Dataset Analysis

Your Name

October 10, 2024

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Dataset Description</b>	<b>2</b>
<b>3</b>	<b>Data Loading and Exploration</b>	<b>2</b>
3.1	Importing Libraries . . . . .	2
3.2	Loading the Dataset . . . . .	3
3.3	Data Overview . . . . .	3
<b>4</b>	<b>Data Cleaning and Preprocessing</b>	<b>3</b>
4.1	Handling Missing Values and Data Types . . . . .	3
4.2	Handling Outliers . . . . .	3
<b>5</b>	<b>Exploratory Data Analysis (EDA)</b>	<b>4</b>
5.1	Distribution of MPG . . . . .	4
5.2	Correlation Matrix . . . . .	4
5.3	Pair Plot . . . . .	5
5.4	Manufacturer Distribution . . . . .	6
5.5	Class Imbalance Analysis . . . . .	7
<b>6</b>	<b>Feature Engineering</b>	<b>8</b>
6.1	Creating New Features . . . . .	8
<b>7</b>	<b>Defining Features and Splitting Data</b>	<b>8</b>
7.1	Preparing Features and Targets . . . . .	8
7.2	Splitting the Data . . . . .	9
<b>8</b>	<b>Modeling and Evaluation</b>	<b>9</b>
8.1	Regression Models . . . . .	9
8.1.1	Linear Regression . . . . .	9
8.1.2	Random Forest Regression . . . . .	10
8.1.3	Support Vector Regression . . . . .	10
8.2	Classification Models . . . . .	10
8.2.1	Logistic Regression . . . . .	10
8.2.2	Random Forest Classification . . . . .	11
8.2.3	Support Vector Classification . . . . .	11
<b>9</b>	<b>Hyperparameter Tuning</b>	<b>12</b>
9.1	Random Forest Regression Tuning . . . . .	12
<b>10</b>	<b>Results Visualization</b>	<b>12</b>
10.1	Regression Results Plot . . . . .	12

# 1 Introduction

In this report, we perform an end-to-end machine learning project using the car dataset. The analysis includes data exploration, preprocessing, feature engineering, modeling, evaluation, and visualization. We aim to predict the miles per gallon (MPG) for cars using regression models and classify whether a car is from Ford using classification models.

## 2 Dataset Description

The **Auto MPG Dataset** consists of 398 cars, including various technical specifications such as:

- **mpg**: Miles per gallon (continuous)
- **cylinders**: Number of cylinders (multi-valued discrete)
- **displacement**: Engine displacement in cubic inches (continuous)
- **horsepower**: Engine horsepower (continuous)
- **weight**: Vehicle weight in pounds (continuous)
- **acceleration**: Time to accelerate from 0 to 60 mph (continuous)
- **model year**: Model year (multi-valued discrete)
- **origin**: Origin of the car (1: USA, 2: Europe, 3: Japan)
- **car name**: Car model name (string)

**Source:** The dataset is available from the UCI Machine Learning Repository.  
<https://archive.ics.uci.edu/ml/datasets/auto+mpg>

## 3 Data Loading and Exploration

### 3.1 Importing Libraries

We start by importing the necessary libraries and ensuring compatibility.

```
1 # Import necessary libraries
2 import sys
3 assert sys.version_info >= (3, 7)
4
5 from packaging import version
6 import sklearn
7 assert version.parse(sklearn.__version__) >= version.parse("1.0.1")
8
9 import pandas as pd
10 import numpy as np
11 import seaborn as sns
12 import matplotlib.pyplot as plt
13 import warnings
14
15 sns.set(style="whitegrid")
16 %matplotlib inline
17
18 warnings.filterwarnings('ignore')
19 np.random.seed(42)
```

## 3.2 Loading the Dataset

We load the dataset using Pandas.

```
1 # Loading the car dataset
2 cars = pd.read_csv("cars.csv")
3
4 # Display the first few rows
5 cars.head()
```

## 3.3 Data Overview

We check data types and identify missing values.

```
1 # Checking data types and missing values
2 cars.info()
3
4 # Summary statistics
5 cars.describe()
```

# 4 Data Cleaning and Preprocessing

## 4.1 Handling Missing Values and Data Types

We handle missing values in the horsepower column and convert data types.

```
1 # Replace '?' with NaN and convert 'horsepower' to numeric
2 cars['horsepower'].replace('?', np.nan, inplace=True)
3 cars['horsepower'] = pd.to_numeric(cars['horsepower'], errors='coerce')
4
5 # Impute missing 'horsepower' with median
6 median_horsepower = cars['horsepower'].median()
7 cars['horsepower'].fillna(median_horsepower, inplace=True)
8
9 # Convert 'origin' to categorical
10 cars['origin'] = cars['origin'].map({1: 'USA', 2: 'Europe', 3: 'Japan'})
11
12 # Extract 'manufacturer' from 'car name'
13 cars['manufacturer'] = cars['car_name'].str.split().str[0]
14
15 # Display cleaned dataset
16 cars.head()
```

## 4.2 Handling Outliers

We detect and remove outliers using the z-score method.

```
1 # Detecting outliers using z-score
2 from scipy import stats
3
4 z_scores = stats.zscore(cars.select_dtypes(include=[np.number]))
5 abs_z_scores = np.abs(z_scores)
6 filtered_entries = (abs_z_scores < 3).all(axis=1)
7 cars = cars[filtered_entries]
8
9 print("Dataset_shape_after_removing_outliers:", cars.shape)
```

## 5 Exploratory Data Analysis (EDA)

### 5.1 Distribution of MPG

We visualize the distribution of the target variable mpg.

```
1 # Distribution of MPG
2 plt.figure(figsize=(10, 6))
3 sns.histplot(cars["mpg"], bins=30, kde=True, color='blue', alpha=0.7)
4 plt.xlabel("Miles_Per_Gallon_(MPG) ")
5 plt.ylabel("Frequency")
6 plt.title("Distribution_of_MPG")
7 plt.savefig("mpg_distribution.png")
8 plt.show()
```

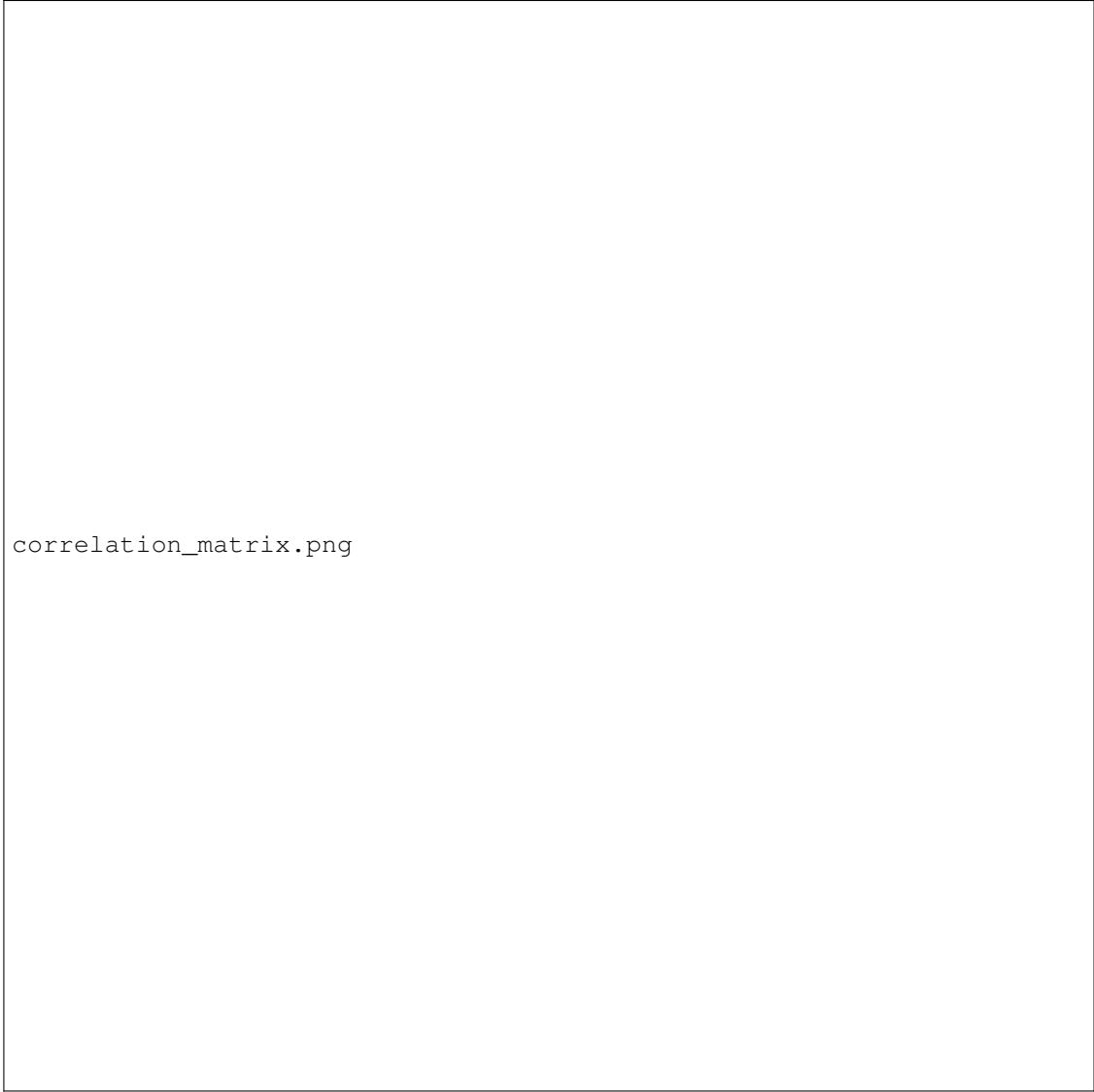


Figure 1: Distribution of MPG

### 5.2 Correlation Matrix

We compute and visualize the correlation matrix.

```
1 # Correlation Matrix
2 plt.figure(figsize=(12, 10))
3 corr_matrix = cars.corr()
4 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm')
5 plt.title('Correlation_Matrix')
6 plt.savefig('correlation_matrix.png')
7 plt.show()
```



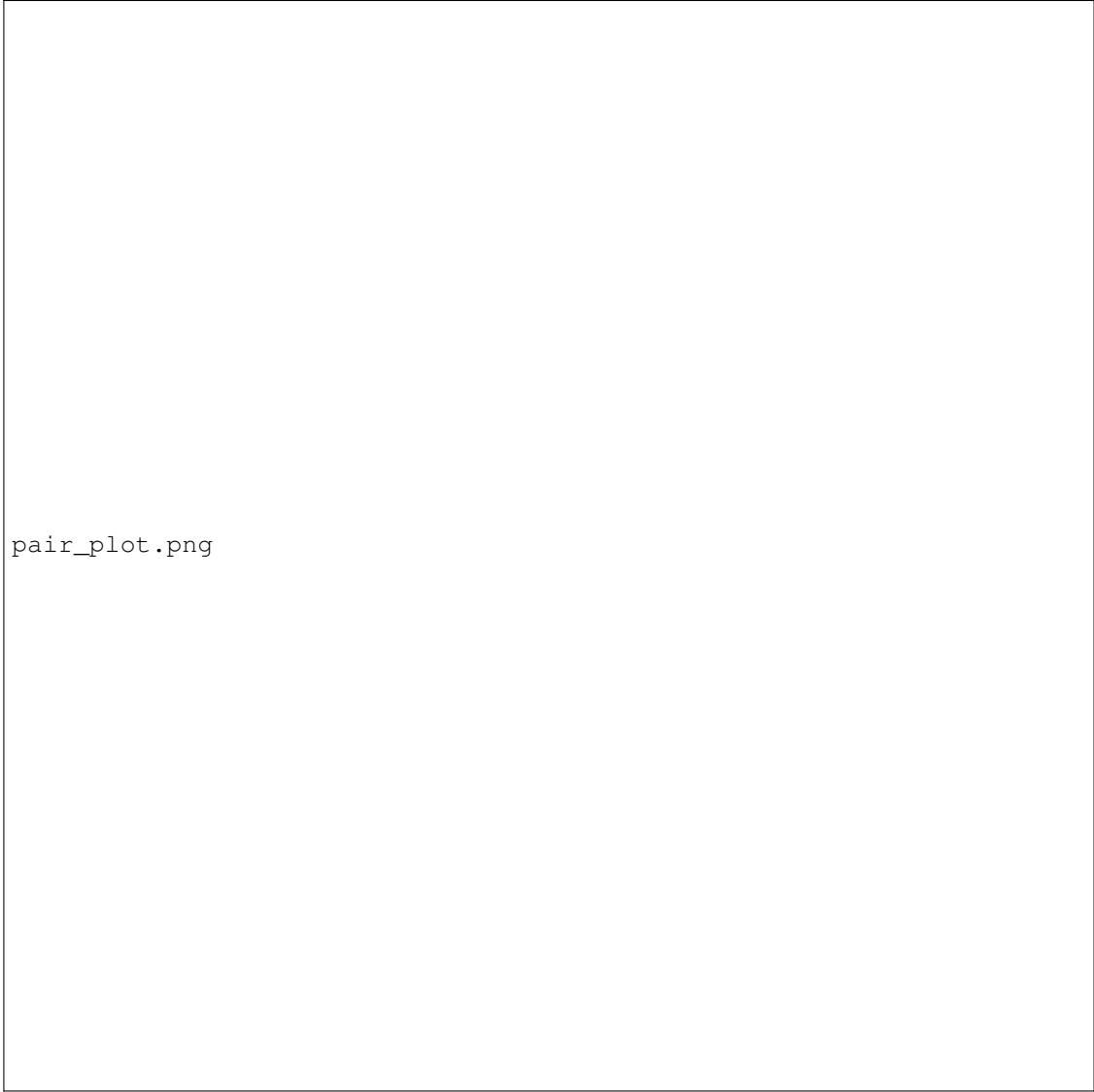
correlation\_matrix.png

Figure 2: Correlation Matrix

### 5.3 Pair Plot

We create pair plots to visualize relationships between features.

```
1 # Pair Plot
2 sns.pairplot(cars[['mpg', 'horsepower', 'weight', 'acceleration', 'displacement']])
3 plt.savefig('pair_plot.png')
4 plt.show()
```



pair\_plot.png

Figure 3: Pair Plot of Selected Features

## 5.4 Manufacturer Distribution

We analyze the distribution of car manufacturers.

```
1 # Distribution of Car Manufacturers
2 plt.figure(figsize=(12, 8))
3 top_manufacturers = cars['manufacturer'].value_counts().nlargest(20)
4 sns.barplot(x=top_manufacturers.index, y=top_manufacturers.values, palette="viridis")
5 plt.xticks(rotation=45)
6 plt.xlabel("Car_Manufacturer")
7 plt.ylabel("Number_of_Cars")
8 plt.title("Top_20_Car_Manufacturers_in_the_Dataset")
9 plt.tight_layout()
10 plt.savefig("manufacturer_distribution.png")
11 plt.show()
```

manufacturer\_distribution.png

Figure 4: Top 20 Car Manufacturers

## 5.5 Class Imbalance Analysis

We examine the class distribution for the classification task.

```
1 # Class Imbalance Analysis
2 cars['is_ford'] = (cars['manufacturer'] == 'ford').astype(int)
3 class_counts = cars['is_ford'].value_counts()
4 print("Class_Distribution:")
5 print(class_counts)
6
7 plt.figure(figsize=(6, 4))
8 sns.barplot(x=class_counts.index, y=class_counts.values)
9 plt.xticks([0, 1], ['Not_Ford', 'Ford'])
10 plt.xlabel('Car_Type')
11 plt.ylabel('Count')
12 plt.title('Ford_vs._Not_Ford_Distribution')
13 plt.savefig('class_distribution.png')
14 plt.show()
```



Figure 5: Ford vs. Not Ford Distribution

## 6 Feature Engineering

### 6.1 Creating New Features

We engineer new features to enhance model performance.

- **Power-to-Weight Ratio:**

$$\text{Power-to-Weight Ratio} = \frac{\text{Horsepower}}{\text{Weight}}$$

- **Displacement per Cylinder:**

$$\text{Displacement per Cylinder} = \frac{\text{Displacement}}{\text{Cylinders}}$$

```
1 # Creating power-to-weight ratio
2 cars['power_to_weight'] = cars['horsepower'] / cars['weight']
3
4 # Creating displacement per cylinder
5 cars['displacement_per_cylinder'] = cars['displacement'] / cars['cylinders']
6
7 # Encoding categorical variables
8 cars = pd.get_dummies(cars, columns=['origin'], drop_first=True)
9
10 cars.head()
```

## 7 Defining Features and Splitting Data

### 7.1 Preparing Features and Targets

We define the feature set and target variables for regression and classification.



```

1 # Defining features and target for regression
2 features = ['horsepower', 'weight', 'acceleration', 'displacement',
3             'power_to_weight', 'displacement_per_cylinder',
4             'origin_Europe', 'origin_Japan']
5 X_regression = cars[features]
6 y_regression = cars['mpg']
7
8 # Defining features and target for classification
9 X_classification = cars[features]
10 y_classification = cars['is_ford']

```

## 7.2 Splitting the Data

We split the data into training and testing sets.

```

1 # Splitting the data into training and testing sets
2 from sklearn.model_selection import train_test_split
3
4 # Regression data split
5 X_train_reg, X_test_reg, y_train_reg, y_test_reg = train_test_split(
6     X_regression, y_regression, test_size=0.2, random_state=42
7 )
8
9 # Classification data split
10 X_train_clf, X_test_clf, y_train_clf, y_test_clf = train_test_split(
11     X_classification, y_classification, test_size=0.2, random_state=42, stratify=y_classification
12 )

```

# 8 Modeling and Evaluation

## 8.1 Regression Models

### 8.1.1 Linear Regression

We train a Linear Regression model and evaluate its performance.

```

1 # Linear Regression Model
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
4
5 lin_reg = LinearRegression()
6 lin_reg.fit(X_train_reg, y_train_reg)
7
8 # Predictions and Evaluation
9 y_pred_lin_reg = lin_reg.predict(X_test_reg)
10 mse_lin_reg = mean_squared_error(y_test_reg, y_pred_lin_reg)
11 mae_lin_reg = mean_absolute_error(y_test_reg, y_pred_lin_reg)
12 r2_lin_reg = r2_score(y_test_reg, y_pred_lin_reg)
13
14 print("Linear Regression Evaluation:")
15 print(f"MSE: {mse_lin_reg:.2f}")
16 print(f"MAE: {mae_lin_reg:.2f}")
17 print(f"R^2 Score: {r2_lin_reg:.2f}")

```

#### Linear Regression Equations:

The prediction is made using:

$$\hat{y} = X\beta + \epsilon$$

where:

- $X$  is the feature matrix

- $\beta$  are the coefficients
- $\epsilon$  is the error term

### 8.1.2 Random Forest Regression

We train a Random Forest Regression model.

```

1 # Random Forest Regression
2 from sklearn.ensemble import RandomForestRegressor
3
4 rf_reg = RandomForestRegressor(random_state=42)
5 rf_reg.fit(X_train_reg, y_train_reg)
6
7 # Predictions and Evaluation
8 y_pred_rf_reg = rf_reg.predict(X_test_reg)
9 mse_rf_reg = mean_squared_error(y_test_reg, y_pred_rf_reg)
10 mae_rf_reg = mean_absolute_error(y_test_reg, y_pred_rf_reg)
11 r2_rf_reg = r2_score(y_test_reg, y_pred_rf_reg)
12
13 print("Random_Forest_Regression_Evaluation:")
14 print(f"MSE:_{mse_rf_reg:.2f}")
15 print(f"MAE:_{mae_rf_reg:.2f}")
16 print(f"R^2_Score:_{r2_rf_reg:.2f}")

```

**Random Forest Regression** combines multiple decision trees to improve predictive accuracy and control over-fitting.

### 8.1.3 Support Vector Regression

We train a Support Vector Regression (SVR) model.

```

1 # Support Vector Regression
2 from sklearn.svm import SVR
3
4 svr_reg = SVR()
5 svr_reg.fit(X_train_reg, y_train_reg)
6
7 # Predictions and Evaluation
8 y_pred_svr_reg = svr_reg.predict(X_test_reg)
9 mse_svr_reg = mean_squared_error(y_test_reg, y_pred_svr_reg)
10 mae_svr_reg = mean_absolute_error(y_test_reg, y_pred_svr_reg)
11 r2_svr_reg = r2_score(y_test_reg, y_pred_svr_reg)
12
13 print("Support_Vector_Regression_Evaluation:")
14 print(f"MSE:_{mse_svr_reg:.2f}")
15 print(f"MAE:_{mae_svr_reg:.2f}")
16 print(f"R^2_Score:_{r2_svr_reg:.2f}")

```

## 8.2 Classification Models

### 8.2.1 Logistic Regression

We train a Logistic Regression model.

```

1 # Logistic Regression Model
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score,
  classification_report
4
5 log_reg = LogisticRegression(max_iter=1000)
6 log_reg.fit(X_train_clf, y_train_clf)
7
8 # Predictions and Evaluation
9 y_pred_log_reg = log_reg.predict(X_test_clf)

```

```

10 accuracy_log_reg = accuracy_score(y_test_clf, y_pred_log_reg)
11 precision_log_reg = precision_score(y_test_clf, y_pred_log_reg)
12 recall_log_reg = recall_score(y_test_clf, y_pred_log_reg)
13 f1_log_reg = f1_score(y_test_clf, y_pred_log_reg)
14
15 print("Logistic_Regression_Evaluation:")
16 print(f"Accuracy:_{accuracy_log_reg:.2f}")
17 print(f"Precision:_{precision_log_reg:.2f}")
18 print(f"Recall:_{recall_log_reg:.2f}")
19 print(f"F1_Score:_{f1_log_reg:.2f}")
20 print("\nClassification_Report:")
21 print(classification_report(y_test_clf, y_pred_log_reg))

```

### Logistic Regression Equation:

The probability of the positive class is modeled as:

$$p = \frac{1}{1 + e^{-(X\beta)}}$$

## 8.2.2 Random Forest Classification

We train a Random Forest Classifier.

```

1 # Random Forest Classification
2 from sklearn.ensemble import RandomForestClassifier
3
4 rf_clf = RandomForestClassifier(random_state=42)
5 rf_clf.fit(X_train_clf, y_train_clf)
6
7 # Predictions and Evaluation
8 y_pred_rf_clf = rf_clf.predict(X_test_clf)
9 accuracy_rf_clf = accuracy_score(y_test_clf, y_pred_rf_clf)
10 precision_rf_clf = precision_score(y_test_clf, y_pred_rf_clf)
11 recall_rf_clf = recall_score(y_test_clf, y_pred_rf_clf)
12 f1_rf_clf = f1_score(y_test_clf, y_pred_rf_clf)
13
14 print("Random_Forest_Classification_Evaluation:")
15 print(f"Accuracy:_{accuracy_rf_clf:.2f}")
16 print(f"Precision:_{precision_rf_clf:.2f}")
17 print(f"Recall:_{recall_rf_clf:.2f}")
18 print(f"F1_Score:_{f1_rf_clf:.2f}")
19 print("\nClassification_Report:")
20 print(classification_report(y_test_clf, y_pred_rf_clf))

```

## 8.2.3 Support Vector Classification

We train a Support Vector Classifier.

```

1 # Support Vector Classification
2 from sklearn.svm import SVC
3
4 svc_clf = SVC(probability=True)
5 svc_clf.fit(X_train_clf, y_train_clf)
6
7 # Predictions and Evaluation
8 y_pred_svc_clf = svc_clf.predict(X_test_clf)
9 accuracy_svc_clf = accuracy_score(y_test_clf, y_pred_svc_clf)
10 precision_svc_clf = precision_score(y_test_clf, y_pred_svc_clf)
11 recall_svc_clf = recall_score(y_test_clf, y_pred_svc_clf)
12 f1_svc_clf = f1_score(y_test_clf, y_pred_svc_clf)
13
14 print("Support_Vector_Classification_Evaluation:")
15 print(f"Accuracy:_{accuracy_svc_clf:.2f}")
16 print(f"Precision:_{precision_svc_clf:.2f}")
17 print(f"Recall:_{recall_svc_clf:.2f}")
18 print(f"F1_Score:_{f1_svc_clf:.2f}")

```

```
19 print("\nClassification_Report:")
20 print(classification_report(y_test_clf, y_pred_svc_clf))
```

## 9 Hyperparameter Tuning

### 9.1 Random Forest Regression Tuning

We use GridSearchCV to fine-tune the Random Forest Regression model.

```
1 # Hyperparameter Tuning for Random Forest Regression
2 from sklearn.model_selection import GridSearchCV
3
4 param_grid_rf_reg = {
5     'n_estimators': [50, 100, 200],
6     'max_depth': [None, 5, 10],
7     'min_samples_split': [2, 5, 10]
8 }
9
10 grid_search_rf_reg = GridSearchCV(RandomForestRegressor(random_state=42), param_grid_rf_reg, cv
    =5)
11 grid_search_rf_reg.fit(X_train_reg, y_train_reg)
12
13 print("Best_Parameters_for_Random_Forest_Regression:")
14 print(grid_search_rf_reg.best_params_)
15
16 # Evaluating the tuned model
17 y_pred_rf_reg_tuned = grid_search_rf_reg.predict(X_test_reg)
18 mse_rf_reg_tuned = mean_squared_error(y_test_reg, y_pred_rf_reg_tuned)
19 mae_rf_reg_tuned = mean_absolute_error(y_test_reg, y_pred_rf_reg_tuned)
20 r2_rf_reg_tuned = r2_score(y_test_reg, y_pred_rf_reg_tuned)
21
22 print("Tuned_Random_Forest_Regression_Evaluation:")
23 print(f"MSE:_{mse_rf_reg_tuned:.2f}")
24 print(f"MAE:_{mae_rf_reg_tuned:.2f}")
25 print(f"R^2_Score:_{r2_rf_reg_tuned:.2f}")
```

## 10 Results Visualization

### 10.1 Regression Results Plot

We plot the actual vs. predicted MPG values for the tuned Random Forest Regression model.