# INDIAN INSTITUTE OF TECHNOLOGY JODHPUR

## Data and Computational Science



# Automobile Price Prediction

## M22MA007
## POTHULA AKASH

# Problem Statement

"Perform data analysis and predictive modeling on automobile pricing data to understand the factors influencing the prices of automobiles. Explore the dataset to identify key features that affect car prices, build predictive models to estimate car prices based on these features, and evaluate the model's performance."

# Colab_link:

https://colab.research.google.com/drive/1DlZcKFf_gj8J-VZz9Q8KCh5b2NKPmoz-?usp=sharing

# Summary of the Project

**Data understanding and cleaning :** examine the Data types of each column , check for the missing values, check the Summary statistics for numeric columns, deal with missing data, converting all the columns to similar datatype

**Exploratory Data analysis:**Drawing box plots ,scatter plots, find correlation,performed statistical tests like ANOVA, Correlation Analysis(finding pearson correlation analysis).

**Machine Learning Models:**  Applied Linear Regression and Multiple Linear Regression.

**Evaluation:** Mean Squared Error , R-squared error, F-test, T-test

# Data before cleaning

# Data after cleaning

# Dealt with missing values using the following methods

**Replacing Missing Values with Averages:** In several instances, we replaced missing values with the average (mean) value of the respective columns. This approach was used for columns such as "normalized-losses," "bore," "stroke," "horsepower," and "peak-rpm."

```python
# Replace missing values in these columns with their respective averages
columns_to_impute = ["normalized-losses", "bore", "stroke", "horsepower", "peak-rpm"]

for column in columns_to_impute:
    # Calculate the average of the current column
    avg_value = df[column].astype("float").mean(axis=0)

    # Replace missing values in the current column with the calculated average
    df[column].replace(np.nan, avg_value, inplace=True)

# Display the calculated average values
for column in columns_to_impute:
    avg_value = df[column].astype("float").mean(axis=0)
    print(f"Average of {column}: {avg_value}")
```

**Filling Missing Values with Most Frequent Values:** For the "num-of-doors" column, missing values were replaced with the most frequent value (mode) in the column.

**Resetting Index:** After dropping rows with missing values, the index was reset using the reset_index method

 **Dropping Rows with Missing Values:** In the "price" column, rows with missing values were simply dropped using the dropna

```python
[11]  # Replace missing values in the 'num-of-doors' column with the most frequent value ('four')
      df['num-of-doors'].replace(np.nan, "four", inplace=True)

      # Drop rows with missing values in the 'price' column
      df.dropna(subset=["price"], axis=0, inplace=True)

      # Reset the DataFrame index after dropping rows
      df.reset_index(drop=True, inplace=True)
```

# Drawing boxplots and histograms

# Pearson correlation coefficient

The Pearson Correlation Coefficient between wheel-base and price is 0.584641822265508 with a P-value of 8.076488270732885e-20

Correlation Strength: The positive value of 0.585 suggests that as the "wheel-base" increases, the "price" of the automobile tends to increase as well. However, the strength of this relationship is moderate, not extremely strong.

Significance: The p-value associated with the correlation coefficient is very close to zero (8.076e-20), indicating that the observed correlation is statistically significant.

```python
from scipy import stats

numeric_columns = ['wheel-base', 'horsepower', 'length', 'width', 'curb-weight', 'engine-size', 'bore', 'city-mpg', 'highway-mpg']

for column in numeric_columns:
    pearson_coef, p_value = stats.pearsonr(df[column], df['price'])
    print(f"The Pearson Correlation Coefficient between {column} and price is {pearson_coef} with a P-value of {p_value}")
```

# Calculating ANOVA

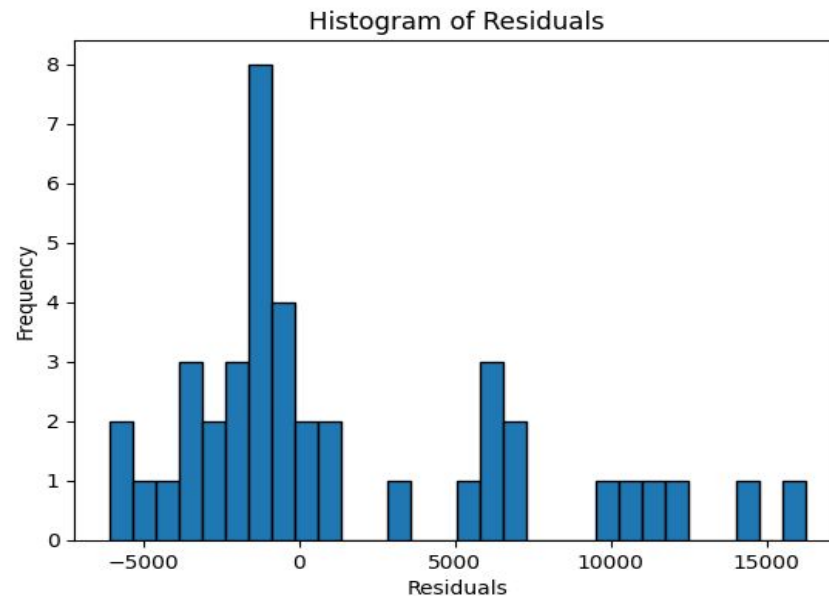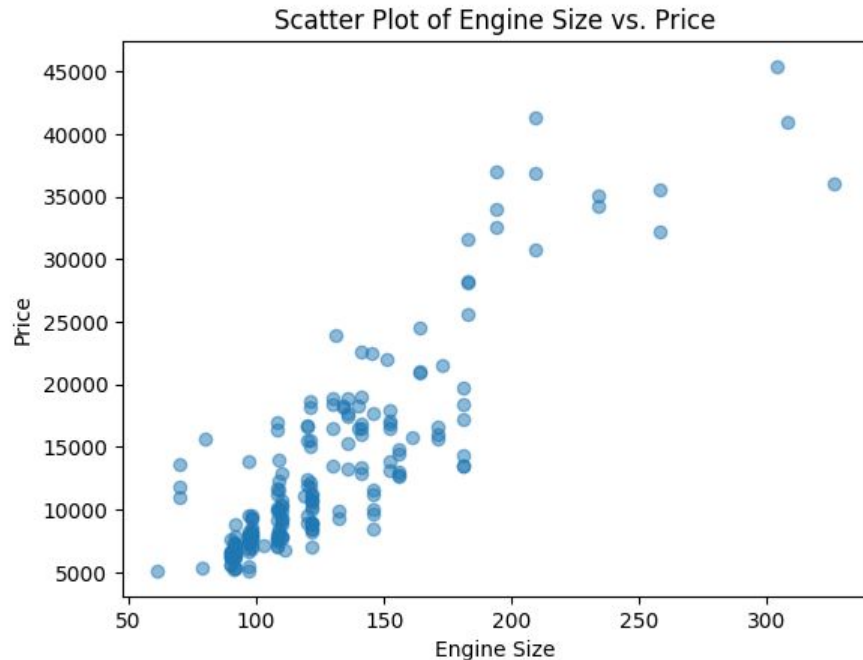ANOVA results for rwd: F=130.5533160959111, P=2.2355306355677845e-23

```python
# List of unique drive-wheels values
drive_wheels = df['drive-wheels'].unique()

# Initialize empty lists to store results
f_values = []
p_values = []

# Calculate ANOVA for each drive-wheels group
for drive_wheel in drive_wheels:
    group = grouped_test2.get_group(drive_wheel)['price']
    f_val, p_val = stats.f_oneway(grouped_test2.get_group('fwd')['price'], group)
    f_values.append(f_val)
    p_values.append(p_val)

# Print ANOVA results
for i, drive_wheel in enumerate(drive_wheels):
    print(f"ANOVA results for {drive_wheel}: F={f_values[i]}, P={p_values[i]}")
```

# Scatter plot and histogram of residuals

# Simple Linear Regression and Q-Q plot

# Multiple Linear Regression

# Evaluation of multiple linear regression and Linear Regression

Mean Squared Error (Multiple Linear Regression): 30393323.64

R-squared (Multiple Linear Regression): 0.75


Mean Squared Error(Linear Regression): 33696986.98421676

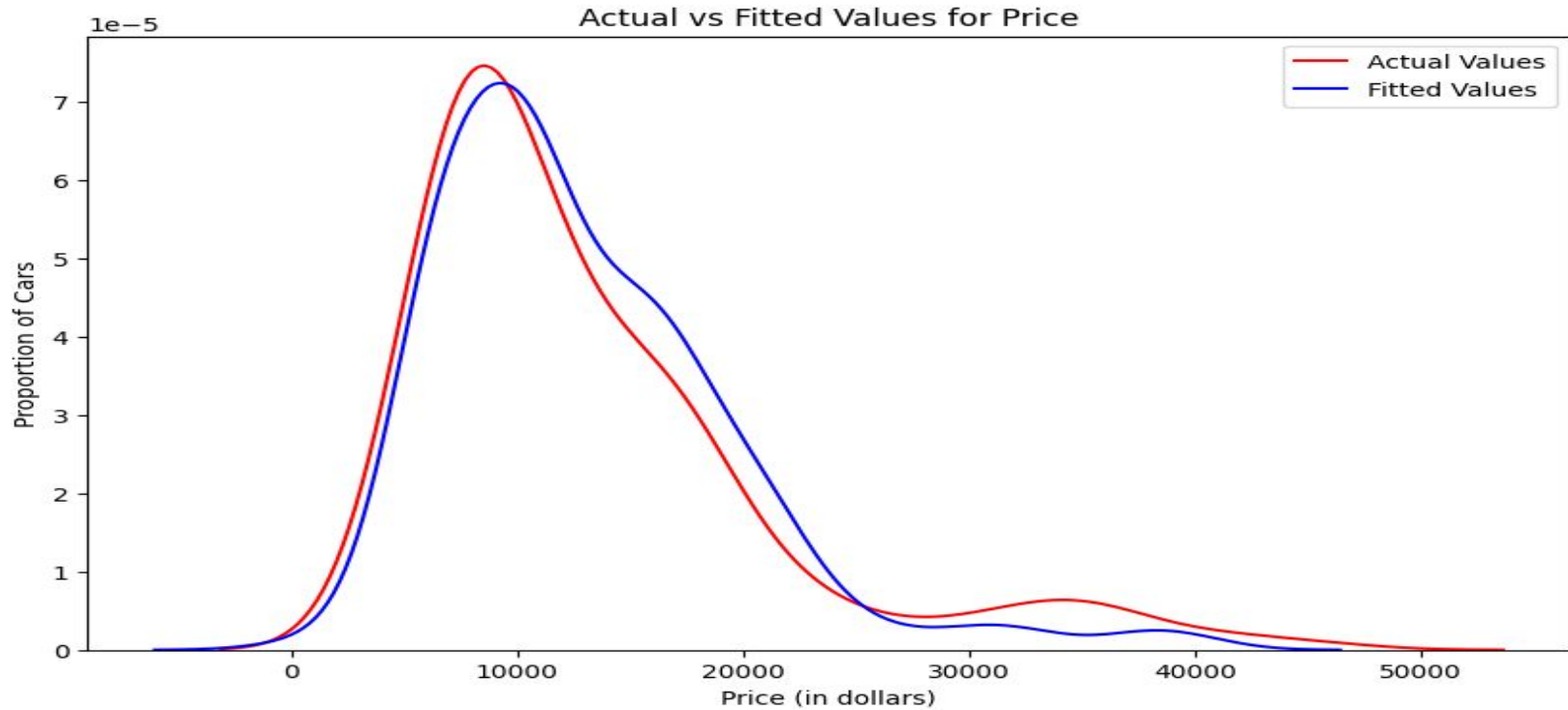R-squared(Multiple Linear Regression): 0.724578048646674

# Overall Model Significance (F-test):

```
Feature: engine-size
F-statistic: 633.5267598010946
P-value: 9.265491622197996e-64
This feature is statistically significant.
==============================
Feature: horsepower
F-statistic: 378.5870228443837
P-value: 6.273536270652618e-48
This feature is statistically significant.
==============================
Feature: curb-weight
F-statistic: 456.138858276953
P-value: 2.189577238897131e-53
This feature is statistically significant.
==============================
Feature: highway-mpg
F-statistic: 356.53919541614164
P-value: 3.0467845810501095e-46
This feature is statistically significant.
==============================
```

# Individual Coefficient Significance (t-test):

```
Feature: const
T-statistic: -10.25497271186749
P-value: 4.9895128140540304e-20
--------------------------
Feature: engine-size
T-statistic: 6.485355428291303
P-value: 7.044275273003609e-10
--------------------------
Feature: horsepower
T-statistic: 2.6067283758896402
P-value: 0.0098438986362382
--------------------------
Feature: curb-weight
T-statistic: 3.28273539858447
P-value: 0.0012172841623101253
--------------------------
Feature: highway-mpg
T-statistic: 1.6674304790368502
P-value: 0.09702580553167357
--------------------------
```

# Actual vs Fitted Values for price

# Thank you