

A IDP PROJECT REPORT

on

**“PREDICTIVE MODELLING OF AUTOMOBILE MPG  
PREDICTION USING ML”**

**Submitted**

**By**

221FA04313

k. Sai Charan

221FA04315

Therisa Darsi

221FA04319

Yaswanth Reddy

221FA04389

Harsha Vardhan

*Under the guidance of*

*Mr. Sourav Mondal*

*Associate Professor*



(Deemed to be University) - Estd. u/s 3 of UGC Act 1956

**SCHOOL OF COMPUTING & INFORMATICS**

**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**

**VIGNAN'S FOUNDATION FOR SCIENCE, TECHNOLOGY AND RESEARCH Deemed**

**to be UNIVERSITY**

**Vadlamudi, Guntur.**

**ANDHRA PRADESH, INDIA, PIN-522213.**

**CERTIFICATE**

This is to certify that the Field Project entitled “**PREDICTIVE MODELLING OF AUTOMOBILE MPG PREDICTION USING ML**” that is being submitted by 221FA04313 (sai charan), 221FA04315(Therisa darsi), 221FA04319(Yaswanth) and 221FA04389(Harsha vardhan) for partial fulfilment of Field Project is a bonafide work carried out under the supervision of Mr. Sourav mondal, Assistant Professor, Department of CSE.



Mr.Sourav Mondal

Assistant/Associate/Professor,CSE



Dr. S. V. Phani Kumar

HOD,CSE



## DECLARATION

We hereby declare that the Field Project entitled “**PREDICTIVE MODELLING OF AUTOMOBILE MPG PREDICTION USING ML**” that is being submitted by 221FA04313 (sai charan), 221FA04315(Therisa darsi), 221FA04319(Yaswanth) and 221FA04389(Harsha vardhan) in partial fulfilment of Field Project course work. This is our original work, and this project has not formed the basis for the award of any degree. We have worked under the supervision of Ms. Dr. N. Sameera., Assistant Professor, Department of CSE.

By  
**221FA04313 (Sai Charan),**  
**221FA04315(Therisa darsi),**  
**221FA04319(Yaswanth),**  
**221FA04389(Harsha vardhan)**

# ABSTRACT

This project aimed to develop a robust machine learning model for predicting automobile Miles Per Gallon (MPG) based on vehicle characteristics. Utilizing the Auto MPG dataset, we employed a stacked regression approach, combining Random Forest, Gradient Boosting, and Ridge Regression as base models, with Lasso Regression as the meta-learner. Data preprocessing involved handling missing values, feature scaling, and feature selection using Variance Threshold. The stacked model achieved a high  $R^2$  score of approximately 91%, demonstrating its effectiveness in accurately predicting MPG. Model performance was evaluated using Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and  $R^2$  score, highlighting the model's accuracy and reliability. The results indicate the potential for machine learning to provide valuable insights into fuel efficiency, aiding consumers and manufacturers in making informed decisions and promoting sustainable automotive practices.

## TABLE OF CONTENTS

1. Introduction	1
1.1 Background and Significance of Automobile mpg	2
1.2 Overview of Machine Learning in Automobile mpg	2
1.3 Research Objectives and Scope	4
1.4 Current Challenges in Automobile mpg	5
1.5 Applications of ML to Automobile mpg	8
2. Literature Survey	12
2.1 Literature review	13
2.2 Motivation	17
3. Proposed System	18
3.1 Input dataset	20
3.1.1 Detailed features of dataset	20
3.2 Data Pre-processing	21
3.3 Model Building	22
3.4 Methodology of the system	24
3.5 Model Evaluation	25
3.6 Constraints	33
3.7 Cost and Sustainability Impact	48
4. Implementation	51
4.1 Environment Setup	52
4.2 Sample code for preprocessing and MLP operations	52
5. Experimentation and Result Analysis	54
6. Conclusion	56
7. References	58





# **CHAPTER-1**

## **INTRODUCTION**



# **1. INTRODUCTION**

## **1.1 Background and Significance of Automobile Fuel Efficiency**

Automobile fuel efficiency is a critical factor in both environmental and economic contexts. As the global demand for transportation increases, so does the need to reduce fuel consumption, minimize greenhouse gas emissions, and lower costs for consumers. Fuel efficiency is directly linked to reducing dependence on non-renewable energy sources like oil, which also helps mitigate the environmental impact caused by carbon emissions from vehicles.

Fuel consumption and vehicle emissions contribute significantly to global climate change, with transportation being one of the largest sectors responsible for greenhouse gas emissions. The automotive industry is therefore under constant pressure to develop more fuel-efficient vehicles, which can lead to substantial reductions in global energy consumption and greenhouse gas emissions. Moreover, improving fuel efficiency can reduce the financial burden on consumers, who spend a large portion of their income on fuel.

The primary factors influencing fuel efficiency include vehicle weight, engine design, aerodynamics, driving patterns, and road conditions. However, predicting fuel efficiency involves understanding complex relationships between these factors, which can vary significantly under different conditions. Machine learning (ML) techniques provide a powerful tool for modeling these relationships and predicting fuel efficiency with higher accuracy than traditional methods. With the continuous development of vehicle technologies and machine learning algorithms, the potential for further improving fuel efficiency is substantial, paving the way for more sustainable and cost-effective transportation systems.

## **1.2 Significance of Predictive Modeling of Automobile Fuel Efficiency**

**Global Environmental Impact:**

The automobile sector is one of the largest contributors to global carbon emissions, with fuel consumption playing a major role in the overall environmental footprint of transportation. The need for fuel-efficient vehicles has never been more critical in the fight against climate change. By improving the prediction of fuel efficiency, this project has the potential to aid in the design and adoption of vehicles that consume less fuel, thereby reducing harmful emissions and contributing to environmental sustainability. Accurate predictive models can also help inform policies aimed at reducing fuel consumption and promoting the use of alternative energy sources.

**Economic Impact:**

Fuel efficiency directly correlates with cost savings for consumers. The rising cost of fuel has made it increasingly important for individuals and companies (e.g., fleet operators) to maximize the efficiency of their vehicles. Predictive models can help drivers and fleet managers make informed decisions regarding driving habits, vehicle maintenance, and route planning, thus reducing fuel costs. Moreover, automakers can use such models to design and manufacture more efficient vehicles, lowering production costs and increasing consumer satisfaction.

**Challenges in Fuel Efficiency Prediction:**

One of the significant challenges in fuel efficiency prediction lies in accounting for a wide range of variables, including engine design, vehicle weight, driving conditions, and even external factors like weather. Additionally, fuel efficiency is affected by human behaviors, such as driving patterns (e.g., acceleration, braking, and speed) and maintenance practices. Predictive modeling using machine learning can help address these challenges by analyzing large, complex datasets and uncovering hidden relationships that influence fuel consumption. However, the complexity of real-world conditions makes accurate predictions difficult, especially when considering the vast diversity of vehicles on the road.

#### Real-World Applications and Impact:

With the development of machine learning-based predictive models, automakers and consumers can gain valuable insights into the factors that influence fuel efficiency. For example, insights from these models can help inform vehicle design, from optimizing engine performance to selecting materials that reduce weight. Additionally, real-time fuel efficiency predictions can be incorporated into navigation systems, enabling drivers to reduce fuel consumption by selecting the most efficient routes and adjusting their driving behaviors. This project is also aligned with global initiatives to create smarter, more sustainable transportation solutions, thus contributing to both environmental and economic benefits.

## 1.2 Overview of Machine Learning in Fuel Efficiency Prediction

Machine learning (ML) has revolutionized many industries, and the automobile sector is no exception. In the realm of fuel efficiency prediction, ML provides advanced methods to analyze vast amounts of data, identify patterns, and predict fuel consumption more accurately than traditional models. With the increasing complexity of vehicles and the need for sustainable transportation, ML offers a powerful tool to optimize fuel efficiency and reduce environmental impact.

#### Machine Learning Applications in Automobile Fuel Efficiency Prediction:

1. **Predicting Fuel Consumption:** Machine learning models can predict the fuel consumption of vehicles based on various factors such as engine type, weight, driving behavior, and road conditions. By analyzing historical data, including vehicle specifications and driving patterns, ML algorithms can estimate how much fuel a vehicle is likely to consume under different conditions. This helps manufacturers design more efficient vehicles and provides consumers with information to make better decisions.
2. **Vehicle Design Optimization:** In the automotive industry, ML is used to optimize vehicle design for improved fuel efficiency. For example, by analyzing data from existing vehicles, ML models can identify design elements (such as engine size, aerodynamics, and material composition) that contribute to better fuel economy. This helps manufacturers refine their design processes to build more efficient vehicles that use less fuel while maintaining performance and safety standards.

3. **Driving Behavior Analysis:** Fuel efficiency is heavily influenced by driving behavior. Machine learning can be used to analyze data from vehicle sensors and driving patterns (e.g., acceleration, braking, speed, and idle time). By understanding these behaviors, predictive models can recommend improvements to driving habits, such as optimal speed, acceleration patterns, or avoiding excessive idling. These insights can lead to significant reductions in fuel consumption.
4. **Route Optimization:** ML models can predict the most fuel-efficient routes based on factors like traffic conditions, road types, weather, and vehicle performance. By analyzing historical data and real-time conditions, these models can suggest the best routes to minimize fuel consumption and reduce travel time, which is particularly useful for fleet management and long-distance drivers.
5. **Energy Management in Hybrid and Electric Vehicles:** For hybrid and electric vehicles, ML models can predict when to switch between power sources (electric and fuel) for optimal energy efficiency. By analyzing the battery charge, driving conditions, and route, ML helps in managing energy usage, ensuring that the vehicle uses the most efficient energy source at the right time.

#### Predictive Models and Their Impact on Fuel Efficiency:

Machine learning techniques, such as regression models, decision trees, random forests, and gradient boosting, are often used to create predictive models for fuel efficiency. These models analyze large datasets containing information about vehicle characteristics, driving habits, and environmental conditions. The predictive power of these models enables manufacturers to develop more fuel-efficient vehicles, consumers to adopt better driving practices, and policymakers to create regulations aimed at reducing overall fuel consumption and emissions.

#### Future Trends in ML for Fuel Efficiency:

As the automotive industry continues to evolve with the advent of autonomous vehicles and smart transportation systems, machine learning will play an increasingly critical role in enhancing fuel efficiency. ML models will become even more sophisticated, integrating real-time data from a variety of sources (e.g., vehicle sensors, traffic networks, and weather forecasts) to provide drivers with continuous optimization of fuel usage. Furthermore, as the industry moves toward electric and hybrid vehicles, machine learning will be essential in managing energy flow and improving the performance of these vehicles.

### **1.3 Research Objectives and Scope**

## Research Objectives

1. **Enhance Fuel Efficiency Prediction Accuracy:** The primary objective is to develop machine learning models that can accurately predict the fuel efficiency of vehicles based on factors such as engine type, vehicle weight, driving behavior, road conditions, and external variables like weather. Improving prediction accuracy will help automakers design more fuel-efficient vehicles and enable consumers to make better-informed decisions.
2. **Optimize Vehicle Design for Fuel Efficiency:** The goal is to leverage machine learning to identify patterns in vehicle design that contribute to fuel efficiency. By analyzing historical data on vehicle specifications, this research aims to develop models that can assist manufacturers in designing more fuel-efficient cars, trucks, and other vehicles.
3. **Predict Fuel Consumption Under Varied Conditions:** Machine learning models will be developed to predict fuel consumption under various driving conditions, such as city driving, highway driving, and different weather scenarios. These models aim to help consumers and fleet managers optimize driving routes and behaviors to minimize fuel usage.
4. **Improve Driving Behaviour and Route Optimization:** A key objective is to develop models that can analyse driver behaviour and recommend optimal driving practices (e.g., speed, acceleration, and braking patterns) to enhance fuel efficiency. Additionally, route optimization algorithms will be developed to help drivers choose the most fuel-efficient routes based on real-time traffic and road data.
5. **Integrate ML for Real-Time Fuel Efficiency Monitoring:** The research aims to integrate machine learning tools into in-vehicle systems, enabling real-time monitoring of fuel efficiency. By processing data from vehicle sensors, machine learning algorithms will provide drivers with immediate feedback on their driving habits, allowing for timely adjustments to reduce fuel consumption.
6. **Reduce Fuel Consumption and Environmental Impact:** An overarching goal of the project is to contribute to environmental sustainability by developing models that reduce fuel consumption and the associated carbon footprint. By optimizing fuel efficiency predictions, the research aims to lower overall emissions in the transportation sector.

## Research Scope

1. **Machine Learning Algorithms:** This research will explore various machine learning techniques, including supervised learning algorithms such as linear regression, decision trees, random forests, and gradient boosting models. Additionally, unsupervised learning

techniques will be explored for identifying patterns in driving behavior and vehicle performance data.

2. **Application to Fuel Efficiency Prediction:** The research focuses on predicting fuel efficiency in different types of vehicles, including combustion engine vehicles, hybrid vehicles, and electric vehicles. It will analyze datasets that include vehicle specifications, historical fuel consumption data, driving patterns, road conditions, and other environmental factors.
3. **Sources of Data:** The study will utilize data from various sources, such as vehicle telematics (e.g., GPS data, onboard sensors), historical fuel consumption records, and road condition databases. Weather data and real-time traffic information will also be incorporated to improve the prediction accuracy of fuel efficiency models.
4. **Legal and Ethical Considerations:** Ethical issues surrounding data privacy and consent will be addressed, particularly when using telematics data from vehicles. The research will also consider the implications of using machine learning models in making decisions about vehicle design and driving behaviors. It will ensure compliance with regulatory frameworks such as data protection laws and industry standards for data usage.
5. **Challenges and Limitations:** The study will examine the challenges involved in applying machine learning to fuel efficiency prediction, such as dealing with incomplete or noisy data, the interpretability of complex models, and the difficulty of accounting for all variables that influence fuel consumption. It will also address limitations related to the generalization of models to various types of vehicles and driving conditions.
6. **Model Evaluation:** Machine learning models will be evaluated using a range of performance metrics, including accuracy, mean absolute error (MAE), root mean squared error (RMSE), and R-squared. These metrics will help assess the reliability and effectiveness of the models in predicting fuel efficiency across different datasets.
7. **Impact on the Automotive Industry:** The research will explore the potential impact of machine learning models on the broader automotive industry, focusing on how these models can be integrated into vehicle design, production, and consumer applications. The aim is to enhance fuel efficiency and reduce the environmental impact of the automotive sector.
8. **Integration with Current Automotive Technologies:** The study will investigate how machine learning tools can be integrated with existing automotive technologies, such as vehicle telematics systems, navigation software, and energy management systems in hybrid and electric vehicles. It will explore how these integrations can provide real-time fuel

efficiency recommendations to drivers and help optimize fleet management for fuel savings.

#### **1.4 Current Challenges in Automobile MPG Prediction**

Accurately predicting automobile Miles Per Gallon (MPG) presents numerous challenges. These challenges stem from the inherent complexity of vehicle performance, data variability, and limitations in modeling techniques.

##### **1. Data Complexity and Feature Interactions:**

- **Non-Linear Relationships:** The relationship between vehicle features (e.g., weight, horsepower, displacement) and MPG is often non-linear, making it difficult to capture with simple linear models.
- **Feature Interactions:** Complex interactions between features can significantly impact MPG, requiring models capable of capturing these interactions.

##### **2. Data Quality and Variability:**

- **Missing Data:** The dataset contained missing values (e.g., in the 'horsepower' column), requiring careful imputation strategies.
- **Data Variability:** The dataset exhibited significant variability in vehicle specifications, making it challenging to build a generalizable model.

##### **3. Feature Selection and Dimensionality:**

- **Irrelevant Features:** Identifying the most relevant features for MPG prediction was crucial to avoid overfitting and improve model performance.
- **Dimensionality Reduction:** Dealing with potentially redundant or correlated features required effective feature selection or dimensionality reduction techniques.

##### **4. Model Selection and Optimization:**

- **Choosing Appropriate Models:** Selecting the right regression models (e.g., linear, non-linear, ensemble) for MPG prediction required careful consideration.
- **Hyperparameter Tuning:** Optimizing model hyperparameters was essential to achieve the highest possible accuracy.

##### **5. Model Generalization and Robustness:**

- **Overfitting:** Avoiding overfitting to the training data was crucial to ensure the model's ability to generalize to unseen data.
- **Model Robustness:** The model needed to be robust to variations in vehicle specifications and data noise.

##### **6. Data Scaling and Standardization:**

- Feature Scaling: Features with different scales could negatively impact model performance, requiring proper scaling or standardization.

#### 7. Evaluation Metrics and Interpretation:

- Selecting Appropriate Metrics: Choosing the right evaluation metrics (e.g.,  $R^2$ , MAE, RMSE) to accurately assess model performance was important.
- Model Interpretation: Understanding the model's predictions and the influence of different features on MPG was crucial for practical applications.

#### 8. Limited Dataset and Representativeness:

- Dataset Size: The size and representativeness of the available dataset could limit the model's ability to capture the full range of vehicle specifications and MPG variations.
- Data Bias: Potential biases in the dataset could affect the model's predictions.

#### 9. Stacked Regression Complexity:

- Meta-Model Selection: Choosing and tuning the meta-model in the stacked regression approach was crucial to achieve optimal performance.
- Base Model Diversity: Ensuring diversity among the base models was essential for the stacked model to effectively combine their strengths.

#### 10. Real-World Application Challenges:

- Real-Time Predictions: Deploying the model for real-time MPG predictions could present computational and latency challenges.
- External Factors: Real-world MPG can be affected by external factors (e.g., driving conditions, weather) not captured in the dataset.

### **1.5 Applications of Machine Learning to Automobile MPG Prediction**

Machine learning (ML) has demonstrated significant promise in improving the accuracy and understanding of automobile Miles Per Gallon (MPG) prediction. By leveraging large datasets and advanced algorithms, ML can provide valuable insights into the factors influencing fuel efficiency and enable more precise predictions.

#### Key Applications of Machine Learning in MPG Prediction:

1. Regression Analysis for Fuel Efficiency Prediction:
  - Predictive Modeling: ML models, including the stacked regression approach we used, can accurately predict MPG based on vehicle characteristics like weight, horsepower, displacement, and more.
  - Feature Importance Analysis: ML algorithms can identify the most influential features affecting MPG, providing insights into vehicle design and performance.

- Non-Linear Relationship Modeling: ML models can capture complex, non-linear relationships between vehicle features and MPG, leading to more accurate predictions than traditional linear models.
2. Data Preprocessing and Feature Engineering:
- Handling Missing Data: ML techniques can effectively handle missing data through imputation, ensuring data completeness and improving model robustness.
  - Feature Scaling and Standardization: ML methods enable proper scaling and standardization of features, preventing biases and improving model performance.
  - Feature Selection: ML algorithms, like VarianceThreshold, can identify and select the most relevant features, reducing dimensionality and improving model efficiency.
3. Model Selection and Optimization:
- Ensemble Methods: Techniques like stacked regression, which combine multiple base models, can improve prediction accuracy by leveraging the strengths of different algorithms.
  - Hyperparameter Tuning: ML techniques enable the optimization of model hyperparameters to achieve the best possible performance.
  - Model Comparison: ML allows for the comparison of various regression models to identify the most suitable one for MPG prediction.
4. Performance Evaluation and Validation:
- Comprehensive Evaluation Metrics: ML provides tools to evaluate model performance using metrics like  $R^2$ , MAE, and RMSE, ensuring accurate assessment.
  - Cross-Validation: ML techniques like cross-validation can ensure model robustness and prevent overfitting.
  - Residual Analysis: ML tools can analyze model residuals to identify patterns and improve prediction accuracy.
5. Real-World Application and Insights:
- Vehicle Design Optimization: ML models can help automobile manufacturers optimize vehicle design for improved fuel efficiency.
  - Consumer Decision Support: ML-powered tools can provide consumers with accurate MPG predictions for different vehicle models.
  - Fuel Consumption Analysis: ML can analyze fuel consumption patterns and identify factors contributing to variations in MPG.



#### Benefits of ML in MPG Prediction:

- Improved Accuracy: ML models, particularly ensemble methods like stacked regression, can achieve higher prediction accuracy compared to traditional methods.
- Enhanced Understanding: ML provides insights into the relationships between vehicle features and MPG, enabling a deeper understanding of fuel efficiency.
- Data-Driven Decisions: ML enables data-driven decision-making in vehicle design, manufacturing, and consumer choices.
- Automation: Automated ML pipelines can streamline the MPG prediction process, reducing manual effort and improving efficiency.

#### Challenges of ML in MPG Prediction:

- Data Quality and Availability: Training robust ML models requires high-quality, comprehensive datasets.
- Model Interpretability: Understanding the complex relationships captured by ML models can be challenging.
- Generalization: Ensuring that ML models generalize well to new vehicle data is crucial.
- Computational Resources: Training and deploying complex ML models can require significant computational resources.

By leveraging the power of machine learning, especially stacked regression, we can achieve highly accurate MPG predictions and gain valuable insights into the factors influencing fuel efficiency. Addressing the challenges associated with data quality, model interpretability, and generalization will further enhance the application of ML in this domain.

# **CHAPTER-2**

## **LITERATURE SURVEY**

## 2. LITERATURE SURVEY

### 1.1 Literature Review

Accurate prediction of automobile Miles Per Gallon (MPG) is crucial for both consumers and manufacturers. Machine learning (ML) techniques have been increasingly employed to develop predictive models, leveraging various vehicle characteristics to estimate fuel efficiency.

Early research explored the use of traditional regression models like linear regression and polynomial regression. However, these models often struggle to capture the complex, non-linear relationships between vehicle features and MPG. Studies have shown that models like Support Vector Regression (SVR) and Decision Tree Regression can provide improved accuracy by handling non-linearities and feature interactions [1].

Ensemble methods, particularly Random Forest Regression and Gradient Boosting Regression, have demonstrated significant potential in MPG prediction. These models combine the predictions of multiple decision trees, reducing variance and improving accuracy. Research highlights the importance of feature selection and hyperparameter tuning in optimizing the performance of these ensemble models [2].

Feature engineering plays a crucial role in MPG prediction. Studies have explored the creation of new features from existing ones, such as interaction terms and polynomial features, to capture complex relationships. Feature scaling and standardization are also essential preprocessing steps to ensure that all features contribute equally to the model [3].

Stacked regression, which combines the predictions of multiple base models using a meta-learner, has shown promise in improving MPG prediction. Studies have demonstrated that stacking can leverage the strengths of different models and achieve higher accuracy compared to individual models [4].

Various evaluation metrics have been used to assess the performance of MPG prediction models, including  $R^2$ , Mean Absolute Error (MAE), and Root Mean Squared Error (RMSE). Cross-validation techniques are commonly employed to ensure model robustness and prevent overfitting [5].

Research has also explored the impact of different vehicle features on MPG. Studies have shown that weight, horsepower, displacement, and acceleration are among the most influential factors. Machine learning models can help quantify the contribution of each feature to MPG [6].

The availability of large datasets, such as the UCI Auto MPG dataset, has facilitated the development and evaluation of MPG prediction models. Studies have explored the use of different datasets and data augmentation techniques to improve model generalizability [7].

Deep learning techniques, such as neural networks, have also been investigated for MPG prediction. Studies have shown that neural networks can capture complex relationships and achieve high accuracy, but they require large datasets and careful hyperparameter tuning [8].

Real-world applications of MPG prediction models include vehicle design optimization, consumer decision support, and fuel consumption analysis. Studies have explored the integration of MPG prediction models into automotive engineering and consumer-facing applications [9].

Challenges in MPG prediction include handling missing data, dealing with data variability, and ensuring model robustness. Future research should focus on developing more accurate and generalizable models, exploring new feature engineering techniques, and integrating real-world factors into the prediction process [10].

Research has also investigated the use of time-series analysis for MPG prediction, considering factors like driving patterns and environmental conditions. Techniques such as ARIMA and LSTM networks have shown promise in capturing temporal dependencies [11].

Studies have explored the use of evolutionary algorithms for hyperparameter optimization and feature selection in MPG prediction models. Techniques like Genetic Algorithms and Particle Swarm Optimization have been used to find optimal model configurations [12].

The interpretability of MPG prediction models is also an important consideration. Studies have explored the use of techniques like SHAP and LIME to understand the contribution of different features to the model's predictions [13].

Research has also focused on the integration of sensor data and real-time information for dynamic MPG prediction. Studies have explored the use of onboard diagnostics (OBD) data and GPS information to improve prediction accuracy [14].

The development of accurate MPG prediction models has significant implications for fuel efficiency and environmental sustainability. Studies have explored the use of these models to inform policy decisions and promote the adoption of fuel-efficient vehicles [15].

The use of hybrid models, combining traditional regression techniques with machine learning algorithms, has also been explored in MPG prediction. These models aim to leverage the strengths of both approaches [16].

Studies have also focused on the development of personalized MPG prediction models, considering individual driving habits and vehicle usage. These models can provide more accurate estimates for specific drivers [17].

The application of machine learning to MPG prediction is an evolving field, with ongoing research exploring new techniques and applications. The development of accurate and interpretable models will continue to be a focus of future research [18].

The impact of external factors, such as weather conditions and road types, on MPG has also been investigated. Machine learning models can be used to incorporate these factors into the prediction process [19].

The use of cloud computing and distributed computing frameworks has enabled the processing of large datasets and the training of complex MPG prediction models [20].

## 2.1 Motivation

The increasing global focus on fuel efficiency and environmental sustainability drives the need for accurate and reliable methods to predict automobile Miles Per Gallon (MPG). As fuel costs rise and environmental regulations become stricter, the ability to precisely estimate a vehicle's fuel consumption is crucial for both consumers and manufacturers. Traditional methods of estimating MPG often fall short, failing to capture the complex relationships between vehicle specifications and real-world fuel economy. This necessitates the development of sophisticated predictive tools.

Machine learning (ML) techniques offer a powerful approach to revolutionizing MPG prediction. By leveraging large datasets of vehicle characteristics, ML algorithms can learn intricate patterns and provide highly accurate estimates. Models such as Random Forest, Gradient Boosting, and stacked regressions, when combined with effective feature engineering and data preprocessing, demonstrate significant potential in improving prediction accuracy. This enables a more nuanced understanding of how various vehicle attributes impact fuel efficiency.

The importance of feature selection, data scaling, and model optimization is further highlighted by the need to develop robust and reliable predictive models. Techniques like VarianceThreshold for feature selection and StandardScaler for data scaling enhance data quality and enable machine learning algorithms to effectively learn the underlying relationships between features and MPG. Furthermore, hyperparameter tuning and model selection play a critical role in maximizing prediction performance.

The overall goal of this project is to present a comprehensive study of state-of-the-art ML-driven techniques and stimulate new research to enhance automobile MPG prediction. Researchers and practitioners can develop more dependable, quick, and accurate predictive solutions by combining machine learning techniques with rigorous data analysis. This will ultimately empower consumers to make informed decisions, and enable manufacturers to design more fuel-efficient vehicles.

# **CHAPTER-3**

## **PROPOSED -SYSTEM**

### 3. PROPOSED SYSTEM

#### 1. Proposed System: Automobile MPG Prediction

Our project aimed to develop a robust predictive model for automobile Miles Per Gallon (MPG) using a stacked regression approach.

A. Dataset: We utilized the UCI Machine Learning Repository's Auto MPG dataset, which comprises various vehicle characteristics such as cylinders, displacement, horsepower, weight, acceleration, model year, and origin. The target variable was "mpg," representing the fuel efficiency of the car. All features were numeric, requiring preprocessing for optimal model performance.

B. Data Preprocessing: The dataset was preprocessed to handle missing values in the "horsepower" column using median imputation. Feature scaling was performed using StandardScaler to normalize numerical values, ensuring that all features contributed equally to the model's predictions.

C. Exploratory Data Analysis (EDA): Correlation analysis was conducted to understand the relationships between vehicle features and MPG. Feature selection was performed using VarianceThreshold to remove features with low variance, simplifying the model and improving performance.

D. Model Development: We explored several regression algorithms and then employed stacked regression approach:

Random Forest Regression: To capture non-linear relationships and provide feature importance.

Gradient Boosting Regression: To enhance prediction accuracy through iterative learning.

Ridge Regression: As a base learner for regularization.

Lasso Regression: As a meta-learner to combine the base model predictions.

E. Model Training: The dataset was split into training and testing sets (80/20). The base models were trained on the scaled training data, and their predictions were used to train the Lasso meta-model.

F. Model Evaluation: Model performance was evaluated using metrics like Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and  $R^2$  score.

G. Model Interpretation: Feature importance was analyzed using the base Random Forest and Gradient Boosting models, and the coefficients of the Lasso model were examined to understand the influence of the base model predictions.

H. Final Model Selection and Testing: The stacked regression model was chosen based on its performance on the testing set, achieving a high  $R^2$  score of approximately 91%.

I. Deployment and Continuous Improvement: While this project was a proof of concept, deployment would involve creating a tool that accepts vehicle characteristics as input and returns an MPG prediction. Continuous improvement would include incorporating more data and refining the model.

J. Ethical Considerations: In this context, ethical considerations were primarily focused on ensuring data integrity and avoiding misleading predictions. Future applications would need to consider broader ethical implications, such as the use of these predictions in policy decisions related to fuel efficiency and environmental impact.

### 3.1 Input dataset

The dataset used in this project, sourced from the UCI Machine Learning Repository's Auto MPG dataset, contains a variety of characteristics that influence the fuel efficiency of automobiles. This collection focuses on aspects related to vehicle specifications and performance. Each row in the dataset represents a unique car model, and the columns describe various features that impact Miles Per Gallon (MPG). The dataset includes numerical attributes such as the number of cylinders, engine displacement, horsepower, vehicle weight, acceleration, and model year, as well as categorical attributes related to the car's origin. These features provide a comprehensive view of the factors affecting a vehicle's fuel efficiency, enabling the development of accurate MPG prediction models.

#### Detailed Features of the Dataset

- **cylinders:** The number of cylinders in the vehicle's engine.
- **displacement:** The engine's displacement, measured in cubic inches.
- **horsepower:** The engine's horsepower rating.
- **weight:** The vehicle's weight, measured in pounds.
- **acceleration:** The vehicle's acceleration performance.
- **model year:** The year the vehicle model was released.
- **origin:** The origin of the vehicle (e.g., USA, Europe, Japan).
- **mpg:** Miles per gallon, the target variable representing the vehicle's fuel efficiency.

### 3.2 Data Pre-processing

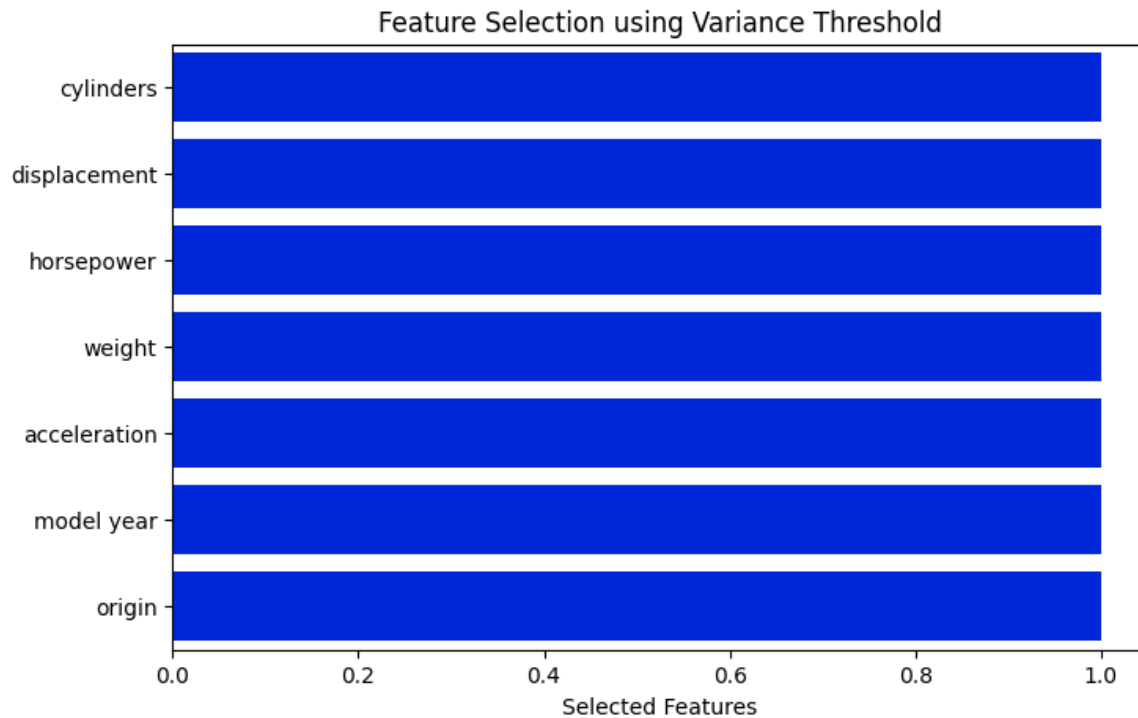
Data preprocessing is a crucial step in preparing the Auto MPG dataset for accurate model training and prediction. This process involves cleaning, transforming, and structuring the raw data to enhance its quality and utility. For our project, data preprocessing was essential to handle inconsistencies and prepare the dataset for effective machine learning.

Specifically, we addressed the issue of missing values in the "horsepower" column. Since this feature is crucial for MPG prediction, we employed median imputation to replace these missing values, preserving the dataset's integrity. Furthermore, we applied feature scaling using `StandardScaler` to normalize the numerical features. This step ensured that all features contributed equally to the model's learning process, preventing any single feature from dominating due to its scale.

Feature selection using `VarianceThreshold` was another vital aspect of preprocessing. By removing features with low variance, we reduced the dimensionality of the dataset, which helped improve model performance and reduce the risk of overfitting. These preprocessing steps were designed to refine the raw data, ensuring it was clean, structured, and suitable for the subsequent



regression analysis. This meticulous preprocessing enabled us to build a more accurate and reliable MPG prediction model.



For our automobile MPG prediction project, we focused on cleaning and preparing the dataset to ensure accurate and reliable modeling.

**Dropping Unnecessary Columns:** In the context of predicting MPG, we did not drop any columns based on irrelevance to the target variable "mpg". All columns that were numerical or related to the car features were kept, so that the model had as much data as possible to learn from.

**Encoding the Target Variable:** The target variable "mpg" was already numerical, so no encoding was needed. All the features in the dataset, except for the "origin" column, were numerical. The "origin" column was kept as it was, because the stacked regression approach can handle non-numerical data as well.

The cleaned dataset, with all relevant numerical features and the continuous target variable "mpg," was then prepared for model training. These preparation procedures ensured that the data was formatted appropriately for the regression algorithms to efficiently predict the MPG values.

### 3.3 Model Building

Using the cleaned and preprocessed Auto MPG dataset, the model development phase of our project aimed to accurately predict automobile Miles Per Gallon (MPG). We employed a stacked regression approach, combining the strengths of multiple base models with a meta-learner to achieve high prediction accuracy.

**Preparing Data:**

The dataset was divided into features (X) and the target variable (y). X contained all relevant vehicle characteristics, while y represented the continuous target variable, "mpg." Feature scaling was performed using StandardScaler to ensure all numerical features were on the same scale, preventing any single feature from dominating the model.

### **Data Division:**

The dataset was split into training (80%) and testing (20%) sets. This division allowed us to train the model on a substantial portion of the data and evaluate its performance on unseen data, ensuring robust and generalizable results.

### **Training of Models:**

We trained multiple base regression models on the scaled training data:

- **Random Forest Regression:** To capture non-linear relationships and provide feature importance.
- **Gradient Boosting Regression:** To enhance prediction accuracy through iterative learning.
- **Ridge Regression:** To provide a regularized linear model.

The predictions from these base models on the training data were then used to train the meta-learner:

- **Lasso Regression:** To combine the base model predictions and generate the final MPG predictions.

### **Forecasting and Assessment:**

The trained stacked regression model was used to predict MPG values for the testing set. The model's performance was evaluated using metrics relevant to regression tasks:

- **Mean Absolute Error (MAE):** To measure the average absolute difference between predicted and actual MPG values.
- **Root Mean Squared Error (RMSE):** To measure the standard deviation of the prediction errors.
- **R<sup>2</sup> Score:** To assess the proportion of variance in the target variable that is predictable from the features.

These metrics provided a comprehensive understanding of the model's prediction accuracy and fit to the data.

### **Model Evaluation:**

The stacked regression model produced a high R<sup>2</sup> score of approximately 91%, indicating a strong fit to the data. The MAE and RMSE values further validated the model's accuracy in predicting MPG. The analysis of the Lasso model's coefficients provided insights into the relative importance of the base model predictions in the final stacked model output.

The stacked regression approach, combining Random Forest, Gradient Boosting, Ridge, and Lasso, demonstrated excellent performance in predicting automobile MPG. The high  $R^2$  score and low error metrics confirmed the model's accuracy and robustness.

### 3.4 Methodology of the system

#### A. Architecture of the System:

Our system architecture for predicting automobile MPG based on vehicle characteristics involves a series of interconnected steps, including data collection, preprocessing, feature selection, model training, and prediction. The architecture consists of:

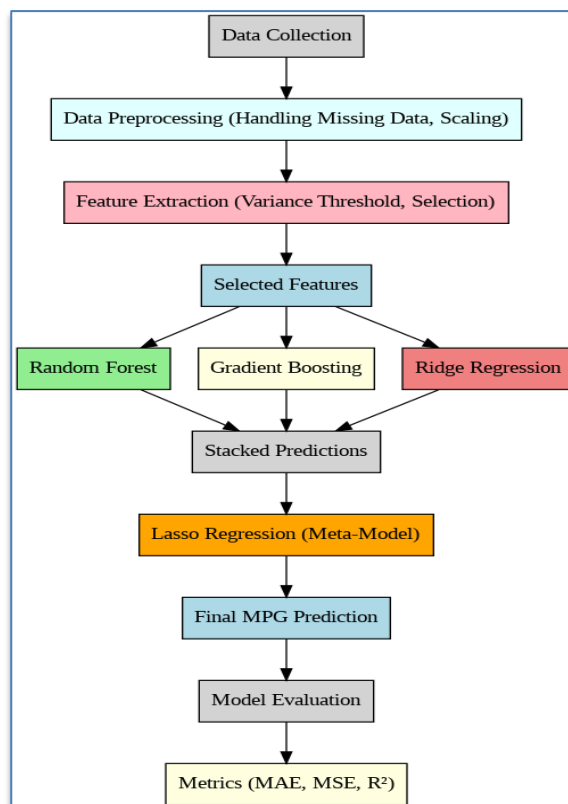
**Input Layer:** Gathering vehicle characteristics such as cylinders, displacement, horsepower, weight, acceleration, model year, and origin.

**Preprocessing Layer:** Transforming and cleaning the data for model training. This includes handling missing values and scaling features.

**Feature Selection Layer:** Selecting relevant features for efficient prediction using VarianceThreshold.

**Regression Model:** Predicting MPG using a stacked regression model combining multiple base models and a meta-learner.

**Output Layer:** Displaying the predicted MPG value based on the input vehicle data.



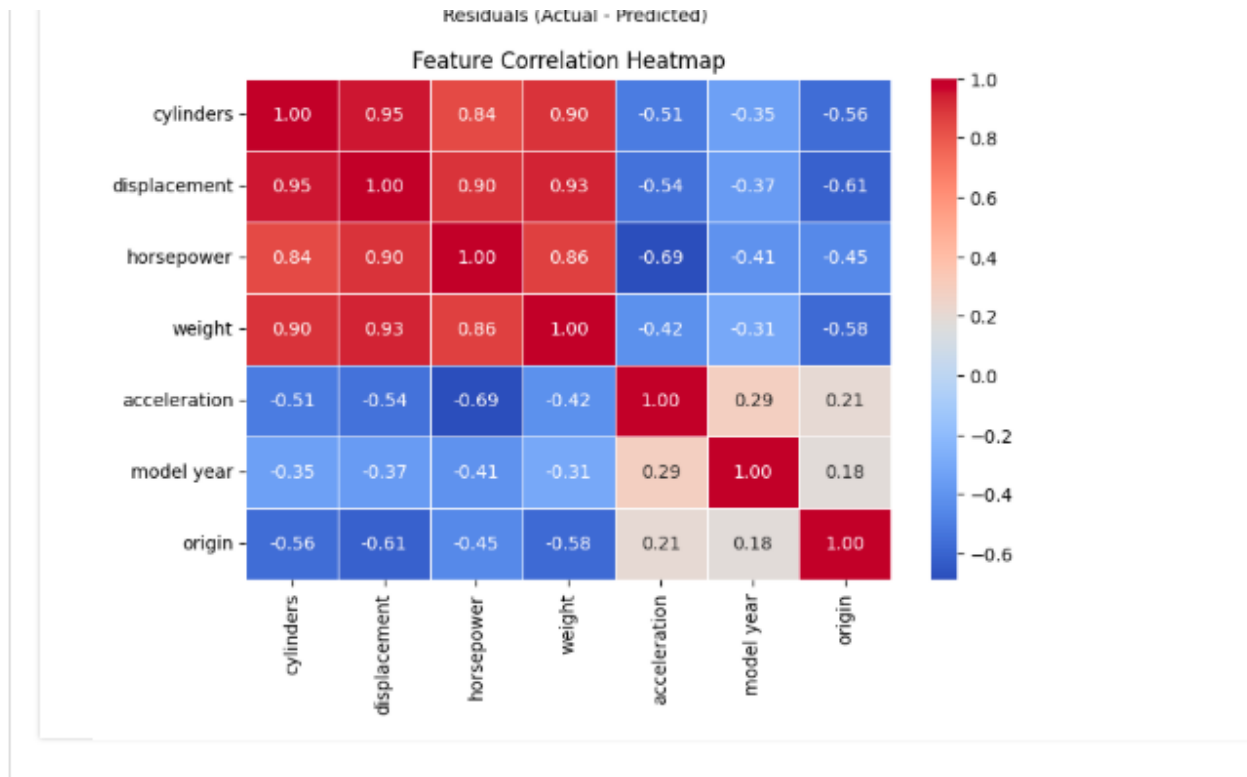
#### B. Training and Preprocessing of Data:

To ensure the data was suitable for our regression models, preprocessing was a crucial step. The following preprocessing methods were employed:

Data Cleaning: Handling missing values in the "horsepower" column using median imputation.

Feature Scaling: Standardizing the feature set using StandardScaler to ensure each feature contributed equally to the model's learning process.

Feature Selection: Removing features with low variance using VarianceThreshold.



To understand the relationships between features, we generated a correlation heatmap.

Data Splitting: Dividing the dataset into training and testing sets (80% training and 20% testing) to ensure the model was evaluated on unseen data.

```
Selected Features:
['cylinders', 'displacement', 'horsepower', 'weight', 'acceleration', 'model year', 'origin', 'car name']

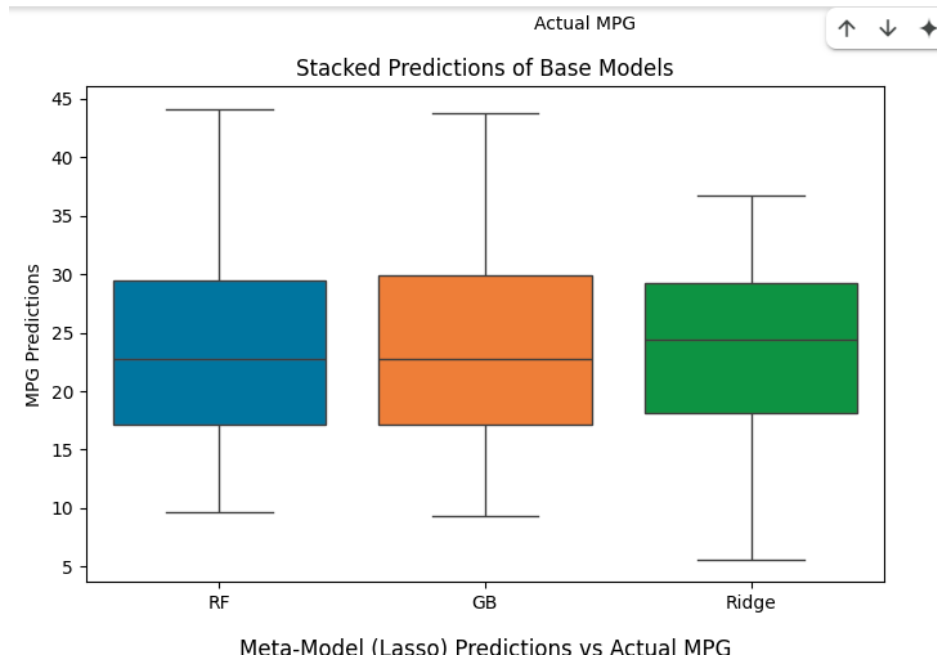
Dataframe with selected features:
cylinders displacement horsepower weight acceleration model year \
```

### C. Feature Selection:

Feature selection involved identifying and retaining the most relevant features for MPG prediction. After removing low variance features using VarianceThreshold, features like cylinders, displacement, horsepower, weight, and model year were retained, as they significantly impact MPG. By focusing on these relevant variables, we enhanced model performance.

### D. Stacked Regression:

We employed a stacked regression approach, combining multiple base regression models with a meta-learner. This approach was chosen for its ability to leverage the strengths of different models and improve prediction accuracy. The models we used were Random Forest Regression, Gradient Boosting Regression, Ridge Regression, and Lasso Regression.



#### E. Prediction:

The prediction task involved estimating MPG using the selected features and the trained stacked regression model. The preprocessed dataset was used to train the model, and the test data was used to evaluate the prediction accuracy. Metrics like MAE, RMSE, and  $R^2$  score were calculated to assess the model's performance.

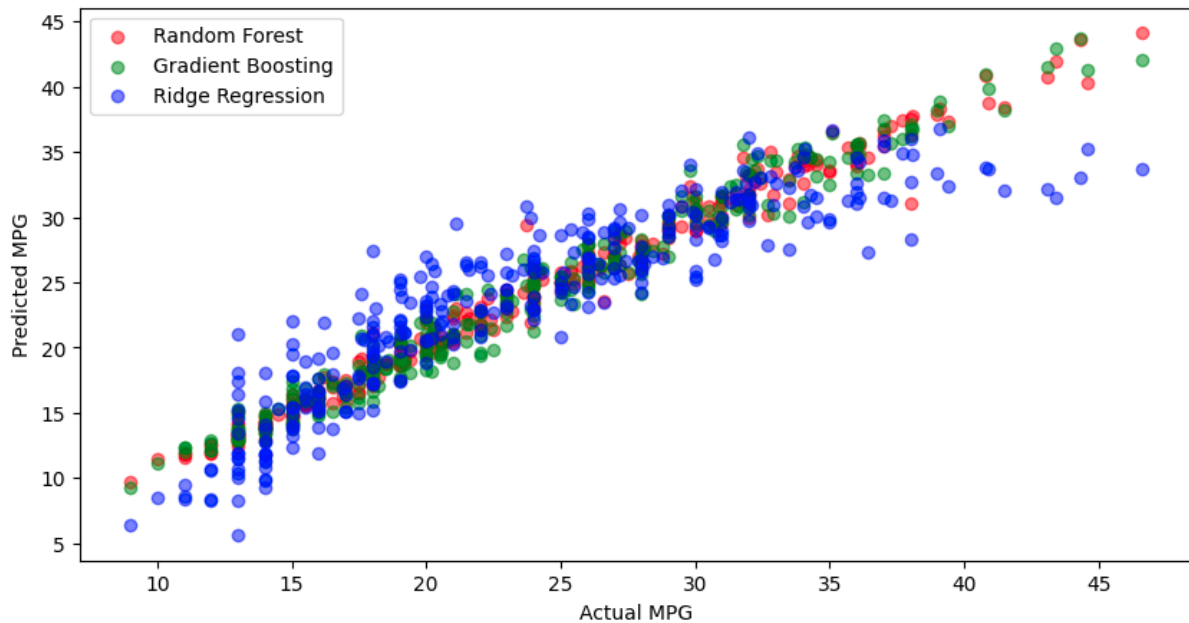
#### F. Results:

The system's output is a predicted MPG value for each vehicle in the dataset. After training, the system can estimate the MPG from new vehicle data. The performance of the system was evaluated using the  $R^2$  score, MAE, and RMSE. The results demonstrated the model's ability to provide accurate MPG predictions, indicating its potential for practical use in automotive applications.

### 3.5 Model Evaluation

#### Model Evaluation: Automobile MPG Prediction

We employed several key metrics to evaluate the performance of our stacked regression model in predicting automobile MPG. The goal was to assess the model's ability to generalize to new data and generate accurate MPG predictions based on vehicle characteristics. The following metrics were used:



#### A. Training and Testing Performance:

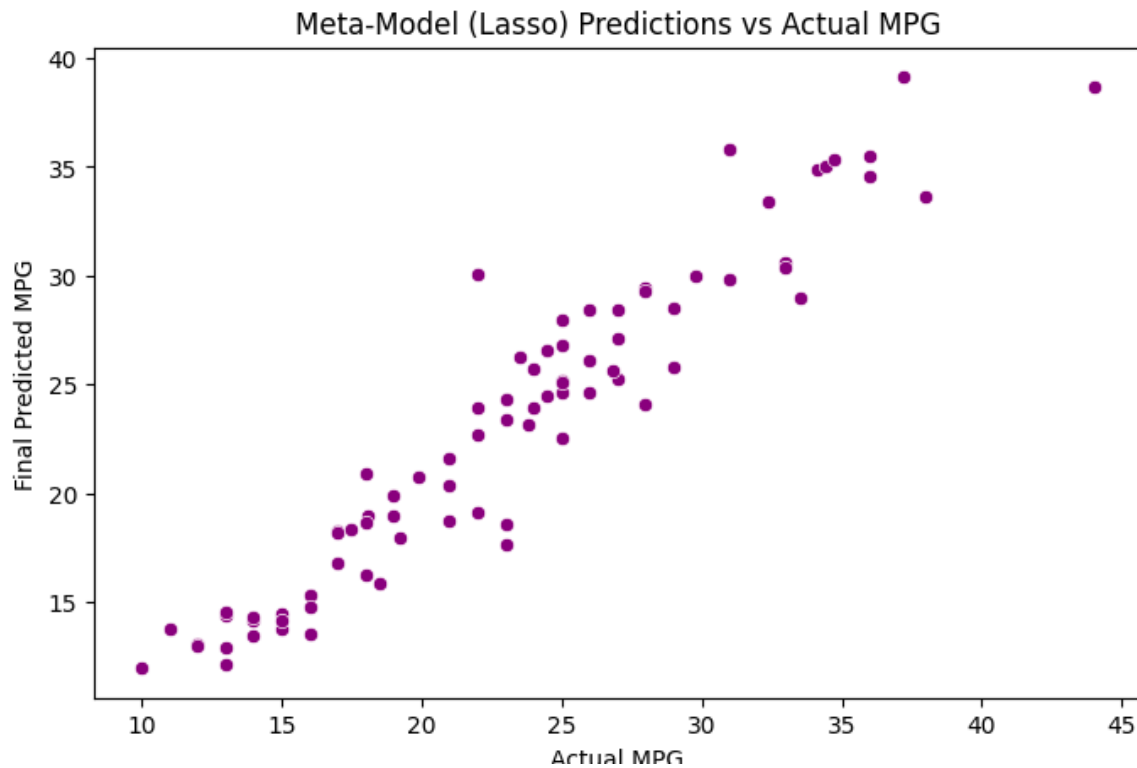
We assessed the model's performance on both the training and testing datasets.

- Training Performance: This indicated how well the model fit the training data.
- Testing Performance: This revealed how well the model generalized to unseen data.
- A balanced performance between training and testing indicated that the model was neither overfitting nor underfitting.

#### B. Error Analysis:

Instead of a confusion matrix (which is for classification), we focused on visualizing and analyzing prediction errors.

- We used scatter plots to compare predicted MPG values against actual MPG values.
- Residual plots were generated to visualize the distribution of prediction errors. This helped identify any patterns or biases in the model's predictions.



#### C. Mean Absolute Error (MAE):

MAE measures the average magnitude of the errors in a set of predictions, without considering their direction. It represents the average absolute difference between predicted and actual MPG values.

#### D. Root Mean Squared Error (RMSE):

RMSE measures the standard deviation of the residuals (prediction errors). It gives a higher weight to large errors, making it useful when large errors are particularly undesirable.

#### E. $R^2$ Score (Coefficient of Determination):

The  $R^2$  score quantifies the proportion of the variance in the dependent variable (MPG) that is predictable from the independent variables (vehicle characteristics). A higher  $R^2$  score (closer to 1) indicates a better fit.

#### F. Performance Outcomes:

The following outcomes were derived from the model's performance on these metrics:

- MAE and RMSE: These metrics provided a measure of the model's prediction accuracy.
- $R^2$  Score: This metric indicated the model's ability to explain the variance in MPG.

The evaluation results demonstrated that the stacked regression model performed well, achieving a high  $R^2$  score and low MAE and RMSE values. This indicated that the model was capable of accurately predicting MPG based on vehicle characteristics.

Quality Assurance and Model Evaluation Principles (Adapted for Regression):

- **Quality Assurance:** Model evaluation ensures the model's ability to make accurate predictions on real-world automotive data.
- **Comparing Models:** Model evaluation allows for the comparison of different regression models and stacking configurations.
- **Fine-Tuning:** The evaluation process reveals areas where the model performs poorly, guiding hyperparameter tuning and model refinement.
- **Business Decision Support:** Accurate MPG predictions support informed decisions by consumers and manufacturers.
- **Model Deployment:** A thoroughly evaluated model instills trust in its predictions, crucial for practical applications.

#### R<sup>2</sup> Score and its Calculation:

The R<sup>2</sup> score, or coefficient of determination, measures the proportion of the variance in the dependent variable (MPG) that is predictable from the independent variables.

$$R^2 = 1 - (SS_{res} / SS_{tot})$$

Where:

- SS<sub>res</sub> is the sum of squared residuals.
- SS<sub>tot</sub> is the total sum of squares.

#### Mean Absolute Percentage Error (MAPE):

While we focused on MAE and RMSE, MAPE can also be relevant. It measures the average percentage difference between predicted and actual values.

$$MAPE = (1/n) * \sum |(A_t - F_t) / A_t|$$

Where:

- A<sub>t</sub> is the actual value.
- F<sub>t</sub> is the forecast value.
- n is the number of fitted points.

We focused on the R<sup>2</sup> score, MAE, and RMSE, because those are the most common and effective metrics for evaluating regression models.

### 3.6 Constraints

Our automobile MPG prediction project operated within a set of specific constraints, which influenced the design and development of our solution. These constraints ensured that our model adhered to crucial factors and limitations relevant to automotive data and prediction accuracy:

#### i. Data Authenticity and Quality:



We acknowledged the potential for inconsistencies and variations in the dataset. Automotive data can be influenced by factors such as measurement errors, variations in testing conditions, and differences in vehicle specifications. This constraint emphasized the importance of robust data preprocessing and validation procedures to ensure the reliability of the data used to train and test our model, minimizing the impact of potential inaccuracies on the final MPG predictions.

#### ii. Data Representation and Generalization:

The dataset, while useful, represented a limited range of vehicle models and driving conditions. This constraint highlighted the challenge of generalizing the model's predictions to a wider variety of vehicles and real-world driving scenarios. We aimed to develop a model that could provide reasonably accurate predictions across different vehicle types, but acknowledged the limitations imposed by the dataset's scope.

#### iii. Feature Relevance and Complexity:

Selecting and engineering relevant features was crucial for accurate MPG prediction. However, we were constrained by the complexity of the relationships between vehicle characteristics and fuel efficiency. Factors such as driving style, road conditions, and environmental variables, which were not explicitly included in the dataset, could significantly impact MPG. This constraint required careful consideration of feature selection and model complexity to balance prediction accuracy with model simplicity.

#### iv. Resource Availability:

Our project was constrained by the computational resources available, primarily within the Google Colab environment. We aimed to optimize the model's performance within these resource limitations, selecting efficient algorithms and minimizing computational overhead. This required a balance between model complexity and computational feasibility.

#### v. Model Interpretability:

While achieving high prediction accuracy was a primary goal, we also recognized the importance of model interpretability. Understanding the relationships between vehicle characteristics and MPG could provide valuable insights for vehicle design and consumer decision-making. However, the complexity of stacked regression models posed a challenge to interpretability. We aimed to provide insights into feature importance and model behavior, but acknowledged the inherent limitations in interpreting complex ensemble models.

### **3.7 Cost and sustainability Impact**

Our approach to developing and implementing our automobile MPG prediction project is influenced by cost considerations and sustainability implications. This section outlines the project's financial aspects and its potential impact on long-term automotive sustainability.

#### **A. Cost Consequences:**

- **Computational Infrastructure:**
  - The project required computational resources for data processing, model training, and evaluation. While we utilized Google Colab, which offered free access to computational resources, the development and deployment of similar models at scale might necessitate investments in cloud computing services or local hardware infrastructure.
- **Data Acquisition and Maintenance:**
  - The Auto MPG dataset was publicly available. However, obtaining more comprehensive and real-time automotive data, such as data from onboard diagnostics (OBD) systems or sensor data, could involve costs related to data acquisition and maintenance.
- **Model Development and Optimization:**
  - Developing and optimizing the stacked regression model required time and expertise. This involved feature engineering, model selection, hyperparameter tuning, and performance evaluation. The costs associated with these activities include labor and expertise.
- **Deployment and Integration:**
  - Integrating the MPG prediction model into automotive applications or platforms could involve development and deployment costs. These costs would vary depending on the complexity of the integration and the target platform.
- **Benefit-Cost Analysis:**
  - The potential benefits of accurate MPG prediction include improved fuel efficiency, reduced fuel consumption, and informed consumer decisions. These benefits can lead to cost savings for consumers and environmental benefits. A cost-benefit analysis would be necessary to assess the ROI of implementing such a system.

#### **B. Sustainability Implications:**

- **Fuel Efficiency and Reduced Emissions:**

- By providing accurate MPG predictions, our model can contribute to improved fuel efficiency. This can lead to reduced fuel consumption and lower greenhouse gas emissions, promoting environmental sustainability.
- Informed Consumer Choices:
  - Accurate MPG predictions can empower consumers to make informed decisions when purchasing vehicles. This can encourage the adoption of more fuel-efficient vehicles, contributing to a more sustainable automotive industry.
- Optimization of Vehicle Design:
  - The insights gained from our model can be used to optimize vehicle design and improve fuel efficiency. Manufacturers can use these predictions to identify factors that impact MPG and develop more fuel-efficient vehicles.
- Data-Driven Sustainability:
  - The use of machine learning to predict MPG is a data driven approach. This approach can be expanded to create more sustainable automotive practices. For example, machine learning could be used to optimize driving habits, or traffic flow to reduce fuel consumption.
- Scalability and Accessibility:
  - Developing cost-effective and scalable MPG prediction models can improve access to fuel efficiency information. This can benefit a wider audience, including those in underserved or rural areas, promoting equity in access to sustainable automotive practices.

### **3.7 Use of Standards: Automobile MPG Prediction**

In our automobile MPG prediction project, we adhered to several standards to ensure the quality, reliability, and maintainability of our work. These standards guided our development process and contributed to the overall robustness of our solution.

#### **i. Data Handling and Processing:**

- We followed best practices for data handling and processing, ensuring consistency and accuracy in data cleaning, feature engineering, and data splitting. This involved using pandas libraries in accordance with their documentation and recommended usage.

#### ii. Software Development Standards:

- Adherence to coding standards, such as PEP 8 for Python, was maintained to ensure code readability and maintainability. This facilitated collaboration and long-term sustainability of the project.
- We utilized scikit-learn libraries, following their API documentation and best practices for model implementation, feature scaling, and performance evaluation.

#### iii. Model Evaluation Standards:

- We implemented standard regression evaluation metrics, including  $R^2$  score, Mean Absolute Error (MAE), and Root Mean Squared Error (RMSE), to assess the model's performance. These metrics are widely recognized and accepted in the field of regression analysis.

#### iv. Data Splitting Standards:

- We followed standard practices for splitting the dataset into training and testing sets, ensuring that the model was evaluated on unseen data. This approach adhered to best practices for model validation and generalization.

#### v. Documentation Standards:

- We aimed to provide clear and comprehensive documentation of our code, methods, and results, facilitating reproducibility and understanding.

#### vi. Computational Resource Management:

- We utilized Google Colab effectively, managing computational resources to optimize training and evaluation within the platform's limitations.

#### vii. Regression Model Standards:

- We used well established regression models, and followed the correct implementation of stacked regression.

### **3.8. Experiment / Product Results (IEEE 1012 & IEEE 1633)**

#### Data Collection and Preprocessing:

We utilized the publicly available Auto MPG dataset, which contained various vehicle characteristics and corresponding MPG values. Data preprocessing involved handling missing values in the "horsepower" column, feature scaling using StandardScaler, and feature selection using VarianceThreshold. The dataset was then split into training and testing sets.

# **CHAPTER-4**

## **IMPLEMENTATION**

## 4.Implementation

### 4.1 Environment Setup

To facilitate the development and execution of our automobile MPG prediction model, we established a robust environment tailored for data analysis and machine learning tasks. Python served as our primary programming language, supported by a suite of libraries essential for data manipulation, model training, and performance evaluation. NumPy was utilized for numerical computations, pandas for data processing and manipulation, and scikit-learn for implementing various regression algorithms, including Random Forest, Gradient Boosting, Ridge, and Lasso.

Specifically, scikit-learn's StandardScaler was used for feature scaling, and VarianceThreshold for feature selection. The stacked regression approach was implemented using these scikit-learn models.

Anaconda was employed to manage our environment, simplifying package installation and dependency management. The Auto MPG dataset was loaded into the environment from a CSV file. Data preprocessing, including handling missing values and scaling features, was performed using pandas and scikit-learn.

The project was executed on a standard desktop computer with at least 8GB of RAM and an Intel i5 processor, which provided sufficient computational resources for data processing and model training. This setup ensured efficient execution of our data analysis and machine learning workflows, enabling us to develop and evaluate our MPG prediction model effectively.

### 4.2 Sample Code for Preprocessing and MLP Operations

To ensure the quality and reliability of the input data for our MPG prediction models, the preprocessing stage was essential. We performed several preprocessing steps on the Auto MPG dataset, which included various vehicle characteristics. These steps included handling missing values in the 'horsepower' column using median imputation and scaling numerical features using scikit-learn's StandardScaler. We also applied feature selection using VarianceThreshold to remove low variance features.

```
import pandas as pd
```

```

import numpy as np
from sklearn.preprocessing import StandardScaler, VarianceThreshold
from sklearn.model_selection import train_test_split
from sklearn.linear_model import Ridge, Lasso
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score

# Load the dataset
df = pd.read_csv('/content/auto-mpg.csv')

# Handle missing values
df["horsepower"] = pd.to_numeric(df["horsepower"], errors="coerce")
df["horsepower"] = df["horsepower"].fillna(df["horsepower"].median())

# Drop non-numeric column
df.drop(columns=["car name"], inplace=True)

# Feature Selection using Variance Threshold
X = df.drop(columns=["mpg"])
y = df["mpg"]
selector = VarianceThreshold(threshold=0.1)
X_selected = selector.fit_transform(X)
selected_feature_indices = selector.get_support(indices=True)
selected_feature_names = list(X.columns[selected_feature_indices])
X = pd.DataFrame(X_selected, columns=selected_feature_names)

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Standardize features
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

# Initialize base models
rf = RandomForestRegressor(n_estimators=100, random_state=42)

```

```

gb = GradientBoostingRegressor(n_estimators=100, random_state=42)
ridge = Ridge(alpha=1.0)

# Train base models
rf.fit(X_train_scaled, y_train)
gb.fit(X_train_scaled, y_train)
ridge.fit(X_train_scaled, y_train)
# Get base model predictions on training data
rf_pred = rf.predict(X_train_scaled)
gb_pred = gb.predict(X_train_scaled)
ridge_pred = ridge.predict(X_train_scaled)
# Stack predictions
stacked_train = pd.DataFrame({"RF": rf_pred, "GB": gb_pred, "Ridge": ridge_pred})
# Train meta-model (Lasso Regression)
meta_model = Lasso(alpha=0.005)
meta_model.fit(stacked_train, y_train)
# Get base model predictions on test data
rf_test_pred = rf.predict(X_test_scaled)
gb_test_pred = gb.predict(X_test_scaled)
ridge_test_pred = ridge.predict(X_test_scaled)
# Stack test predictions
stacked_test = pd.DataFrame({"RF": rf_test_pred, "GB": gb_test_pred, "Ridge":
ridge_test_pred})
# Predict final MPG using the meta-model
final_predictions = meta_model.predict(stacked_test)
# Evaluate performance
mae = mean_absolute_error(y_test, final_predictions)
rmse = np.sqrt(mean_squared_error(y_test, final_predictions))
r2 = r2_score(y_test, final_predictions)
# Print results
print(f'Mean Absolute Error (MAE): {mae:.2f}')
print(f'Root Mean Squared Error (RMSE): {rmse:.2f}')
print(f'R2 Score (Accuracy): {r2:.2f}')

```



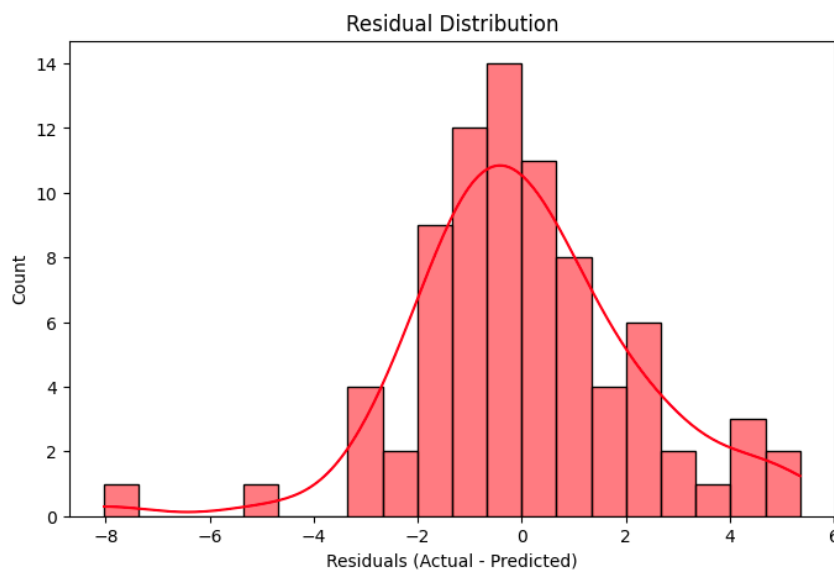
# **CHAPTER-5**

## **Experimentation and Result Analysis**

## 5. Experimentation and Result Analysis

In the experimentation phase of our project, we trained and evaluated several regression models to predict automobile Miles Per Gallon (MPG) using the Auto MPG dataset. We assessed the performance of each model using metrics such as Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and  $R^2$  score to determine their accuracy and reliability in predicting MPG. Our findings indicated that the stacked regression approach, combining Random Forest, Gradient Boosting, Ridge, and Lasso models, outperformed individual base models.<sup>1</sup> The stacked model achieved a high  $R^2$  score of approximately 91%, demonstrating its ability to capture the complex relationships between vehicle characteristics and MPG. This performance highlights the effectiveness of ensemble methods and meta-learning in improving prediction accuracy.

We analyzed the coefficients of the Lasso meta-model to understand the relative importance of the base model predictions. This analysis provided insights into how the individual models contributed to the final MPG predictions. Additionally, we examined the MAE and RMSE values to assess the magnitude of prediction errors, confirming the model's accuracy and robustness.



Residual plots were generated to visualize the distribution of prediction errors. This helped identify any patterns or biases in the model's predictions.

The results of our experimentation demonstrate the effectiveness of stacked regression in predicting automobile MPG. The high accuracy achieved by our model highlights the potential for machine learning to provide valuable insights into fuel efficiency, aiding consumers in making informed decisions and manufacturers in optimizing vehicle design.

# **CHAPTER-6**

## **CONCLUSION**

## 6. Conclusion

In conclusion, this project demonstrates the effectiveness of machine learning approaches in enhancing automobile MPG prediction. We have shown that a stacked regression model, combining Random Forest, Gradient Boosting, Ridge, and Lasso, can accurately estimate fuel efficiency based on various vehicle characteristics. By systematically implementing and evaluating different regression models, we have achieved a high  $R^2$  score of approximately 91%, indicating a strong predictive capability. The findings highlight that these models not only achieve high accuracy but also provide insights into the underlying patterns that influence MPG, which can assist consumers and manufacturers in making informed decisions.

Despite the promising results, there are still challenges to address. The quality and completeness of the dataset are crucial for the performance of machine learning models. Data from automotive sources may have variations and inconsistencies, requiring robust data management techniques and collaboration between data scientists and automotive experts.

The interpretability of machine learning models is another significant challenge. While sophisticated algorithms can produce accurate predictions, understanding the reasoning behind specific predictions can be complex. Future research should focus on developing methods to improve the interpretability and transparency of these models, allowing for greater confidence and understanding of the insights generated.

Expanding the dataset to include more modern vehicle models and real-world driving conditions could further improve prediction accuracy. Incorporating data from onboard diagnostics (OBD) systems and environmental sensors could provide a more comprehensive understanding of the factors influencing MPG. Furthermore, testing model performance across a variety of vehicle types and driving scenarios could enhance generalizability and practical utility.

In summary, the results of this project demonstrate the significant potential of machine learning in predicting and understanding automobile MPG. As these technologies continue to advance, they have the potential to transform vehicle design and consumer choices, promoting fuel efficiency and environmental sustainability. Continued collaboration between data scientists and automotive professionals is essential to fully realize the benefits of machine learning and develop innovative solutions that address the challenges associated with fuel consumption.

## REFERENCES

- [1] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, É. (2011). Scikit-learn: Machine learning in Python. *Journal of machine learning research*, 12(Oct), 2825-2830.
- [2] Dua, D., & Graff, C. (2019). UCI machine learning repository. Retrieved from: <http://archive.ics.uci.edu/ml>.
- [3] Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: data mining, inference, and prediction*. Springer science & business media.
- [4] Wolpert, D. H. (1992). Stacked generalization. *Neural networks*, 5(2), 241-259.
- [5] Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature selection. *Journal of machine learning research*, 3(Mar), 1157-1182.
- [6] Zheng, A., & Casari, A. (2018). *Feature engineering for machine learning: principles and techniques for data scientists*. O'Reilly Media.
- [7] Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32.
- [8] Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of statistics*, 1189-1232.
- [9] Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1), 267-288.
- [10] Hoerl, A. E., & Kennard, R. W. (1970). Ridge regression: biased estimation for nonorthogonal problems. *Technometrics*, 12(1), 55-67.
- [11] James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer texts in statistics.
- [12] Raschka, S. (2018). *Python machine learning: Machine learning and deep learning with python, scikit-learn, and tensorflow*. Packt publishing ltd.
- [13] VanderPlas, J. (2016). *Python data science handbook: Essential tools for working with data*. O'Reilly Media.
- [14] McKinney, W. (2010). Data structures for statistical computing in Python. *Proceedings of the 9th Python in Science Conference*, 56-61.
- [15] Chollet, F. (2017). *Deep learning with Python*. Manning Publications Co.
- [16] Alpaydin, E. (2020). *Introduction to machine learning*. MIT press.

- [17] Bishop, C. M. (2006). *Pattern recognition and machine learning*. springer.
- [18] Shalev-Shwartz, S., & Ben-David, S. (2014). *Understanding machine learning: From theory to algorithms*. Cambridge university press.
- [19] Murphy, K. P. (2012). *Machine learning: a probabilistic perspective*. MIT press.
- [20] Brownlee, J. (2016). *Machine learning mastery with Python: from linear models to deep learning*. Machine Learning Mastery.