



# How to Quantify Automotive Luxury?

**Boyang Wan**

MGSC 661: Multivariate Statistics

Instructor: Juan Camilo Serpa

Dec 3, 2024

# Introduction

The concept of luxury in the automotive industry is often subjective, driven by perception rather than quantitative metrics. **The objective of this study is to quantify the “luxury perception” associated with car brands**, transforming it into a measurable score. Luxury, as an abstract concept, is difficult to define and standardize, especially for car brands. To address this, a **luxury score** was formulated by engineering features that could objectively capture the luxury aspects of a car.

This analysis also aims to understand **market segmentation** by clustering different car models with their luxury scores, thereby offering business insights that can help brands position their models strategically compared to competitors. These insights are geared towards helping brands adjust specific features to either enhance their luxury perception or position themselves as economical alternatives to target different consumer segments.

## Data Exploration and Feature Engineering

### Dataset Overview

The dataset used in this analysis includes various attributes of cars such as make, fuel type, engine size, curb weight, etc., totaling 193 observations across 39 variables. For a part of detailed descriptions of features, refer to Table 1.

### Data Preprocessing

To prepare the dataset for analysis, several preprocessing steps were undertaken. These included handling missing values, standardizing numerical features, encoding categorical variables, and normalizing the data to ensure compatibility between different features. Additionally, feature scaling was performed to ensure that attributes with larger numerical ranges did not disproportionately affect the model training.

### Feature Engineering

To quantify the luxury of car models, multiple new features were engineered:

#### 1. Power-to-Weight Ratio

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

Higher power-to-weight ratios are typically associated with performance and luxury.

#### 2. Size Index

$$\text{Size Index} = \text{Wheel Base} \times \text{Width} \times \text{Height}$$

The overall size of a car may indicate its luxury status (e.g., larger cars are often considered more premium).

#### 3. Luxury Brand Indicator

$$\text{Luxury Brand Indicator} = \begin{cases} 1 & \text{if Make is Jaguar, Mercedes-Benz, Porsche, BMW, or Volvo} \\ 0 & \text{otherwise} \end{cases}$$

These brands were selected based on their high positive price coefficients from the regression analysis, indicating strong associations with higher prices and perceived luxury. Their established reputation in the premium market segment supports their inclusion, ensuring the indicator reflects true luxury brands. For a visualization of the price coefficients by car brands, refer to the figure in the appendix (Figure 1).

#### 4. Fuel Economy Difference

$$\text{Fuel Economy Difference} = \text{Highway MPG} - \text{City MPG}$$

The difference between highway and city mileage can indicate performance tuning for efficiency.

#### 5. Performance Index

$$\text{Performance Index} = \frac{\text{Engine Size} \times \text{Compression Ratio} \times \text{Peak RPM}}{1000}$$

The Performance Index is a composite metric combining these three factors to provide a holistic view of a car's engine performance:

- **Engine Size:** Contributes to the raw power potential.
- **Compression Ratio:** Indicates efficiency and power optimization.
- **Peak RPM:** Reflects the engine's ability to deliver power quickly.

By combining these, the Performance Index captures a balance between raw power, efficiency, and speed capability, which are key dimensions of a car's performance. The division by 1000 is applied to scale down the values for better interpretability and to prevent extremely large numbers from dominating the analysis.

#### 6. Weight-to-Size Ratio

$$\text{Weight-to-Size Ratio} = \frac{\text{Curb Weight}}{\text{Wheel Base} \times \text{Width} \times \text{Height}}$$

This captures the weight distribution relative to the car's size, providing insights into how efficiently the car's weight is managed for its dimensions, which can affect performance and handling.

#### 7. Compression Efficiency

$$\text{Compression Efficiency} = \frac{\text{Compression Ratio} \times \text{Horsepower}}{\text{Engine Size}}$$

This metric reflects how efficiently the engine converts fuel to power by balancing the compression ratio and horsepower with engine size, providing insights into the engine's performance optimization.

#### 8. Doors-to-Weight Ratio

$$\text{Doors-to-Weight Ratio} = \frac{\text{Number of Doors}}{\text{Curb Weight}}$$

This metric reflects accessibility and weight distribution, indicating how the number of doors scales with the vehicle's overall weight for practical and design considerations.

## Multicollinearity and Feature Reduction

During the initial feature selection process, multicollinearity among the variables was addressed by calculating the Variance Inflation Factor (VIF) for each feature. Features with high VIF scores were grouped and dropped based on redundancy and high correlation.

The features dropped are:

- **Highly Correlated with Size Index:**
  - *Wheel Base, Length, Width, Height*: Dropped due to high correlation with *Size Index*, which effectively captured the overall vehicle dimensions.
- **Redundant Weight Metrics:**
  - *Curb Weight, Engine Size*: Removed because they were highly correlated with *Power-to-Weight Ratio, Size Index*, and *Performance Index*. These metrics already represented the necessary relationships between weight, size, and engine characteristics.
- **Redundant Performance Indicators:**
  - *Horsepower, Compression Ratio*: Dropped due to their contribution already being encapsulated within *Power-to-Weight Ratio* and *Performance Index*.
- **Fuel Efficiency Overlap:**
  - *Highway MPG, City MPG*: Dropped as their difference (*Fuel Economy Difference*) provided a more meaningful metric for efficiency.
- **Low Predictive Categorical Variables:**
  - *Make*: Dropped in favor of *Luxury Brand Indicator*, which better represented the luxury status of the car brands.
  - *Number of Doors, Number of Cylinders*: Removed due to limited contribution to luxury perception and redundancy with other metrics.

## Model Selection and Methodology

### Feature Selection

The primary objective of feature selection was to determine the key factors that could influence car prices, which served as a proxy for the concept of "luxury." To achieve this, both linear regression and random forest models were employed to rank and identify the most important features related to price. These selected features would later form the basis of the luxury score.

- **Linear Regression:** Linear regression was initially used to assess the relationship between each feature and the target variable, *price*. The regression coefficients were ranked to determine the most influential features, providing a straightforward approach for feature interpretability.
- **Random Forest:** A Random Forest model was also applied to capture complex non-linear interactions between features. The model provided an importance ranking for each feature, identifying which attributes contributed most significantly to predicting the target variable.

- **Final Feature Selection:** The results from both models were combined and normalized. The top 10 features with the highest average importance scores across both models were selected for further analysis. This approach ensured that only the most impactful features, as confirmed by multiple methods, were used in constructing the luxury score. For a visualization of the top 10 selected features by their normalized importance scores, see Figure 2 in the appendix.

## Luxury Score Construction

The **Luxury Score** for each car is calculated as:

$$Luxury\_Score_i = \sum_{j=1}^N (w_j \cdot x_{ij})$$

Where:

- $i$ : Represents the car (row in the dataset).
- $j$ : Represents the feature (10 most important features selected using linear regression and random forest).
- $N$ : Total number of selected features.
- $w_j$ : The weight of the  $j$ -th feature, derived from the **Average Normalized Importance** across both linear regression and random forest.
- $x_{ij}$ : The normalized value of the  $j$ -th feature for the  $i$ -th car.
- $Luxury\_Score_i$ : The calculated luxury score for the  $i$ -th car.

Each car's luxury score is determined by summing the product of its normalized feature values and their corresponding feature importance weights.

## Detailed Explanation

The luxury score calculation is designed to quantify the perceived luxury level of each car by integrating the contributions of the selected features, weighted by their relative importance. Below, we break down the components of the formula in more detail:

1. **Selected Features ( $x_{ij}$ ):** The features used in this formula are the top 20 features identified during the feature selection process. Each feature represents an important attribute that significantly affects the price, which acts as a proxy for luxury perception.
2. **Normalization of Features:** Each feature value ( $x_{ij}$ ) is normalized to bring all values to a common scale, typically between 0 and 1. This prevents any feature with larger raw values from disproportionately influencing the luxury score.
3. **Weights ( $w_j$ ):** The weights ( $w_j$ ) are derived from the average importance of each feature, as determined by both the linear regression and random forest models. The importance scores from each model were normalized and averaged to ensure a balanced representation of feature significance, accounting for both linear and non-linear effects.

4. **Weighted Sum for Luxury Score** ( $Luxury\_Score_i$ ): The final luxury score for each car is calculated by taking the weighted sum of the normalized feature values. This ensures that each feature contributes to the luxury score proportionally to its importance, providing a composite measure that reflects the combined effect of all selected features on the perception of luxury.

The resulting **Luxury Score** is a numerical value that allows for comparison across different car models, indicating their relative luxury level. Cars with higher luxury scores are those that perform well across the most significant features, thus embodying characteristics that contribute strongly to the perception of luxury. Refer to Appendix Figure 3 and Figure 4 for the distribution of luxury scores and the relationship between luxury scores and car prices, respectively.

## Clustering Analysis

A clustering methodology was employed to segment car models into distinct groups for market positioning insights. The goal was to understand how different car features relate to perceived luxury and identify natural groupings in the data that could guide business decisions.

### Dataset Preparation

A subset of 10 key features was selected for clustering, including engineered features such as:

- `power_to_weight_ratio`
- `size_index`
- `luxury_brand_indicator.1` (representing luxury branding)
- `fuel_economy_difference`
- `performance_index`
- `weight_to_size_ratio`
- `compression_efficiency`
- `doors_to_weight_ratio`
- `Luxury_Score`
- `price`

The inclusion of `Luxury_Score` and `price` was intended to ensure that the clustering captured both the engineered aspects of luxury and the market valuation.

### Methodology

K-means clustering was used to segment the car models into clusters based on the selected features. The clustering process was guided by silhouette analysis to determine the optimal number of clusters. As shown in Figure 5, the silhouette scores for cluster numbers ranging from 2 to 10 were evaluated, and the optimal number of clusters ( $k$ ) was determined to be 3. The resulting clusters and their spatial distribution are visualized in Figure 6.

# Results

## Feature Importance for Luxury Score

To determine which features most significantly contribute to the luxury score, both linear regression and random forest models were used. Features were ranked based on importance from each model, normalized, and then aggregated to select those with the highest significance.

These top features effectively quantify the luxury aspects of car brands, aligning with the objective to make the concept of luxury measurable. The combined rankings from linear and non-linear models capture a holistic view of luxury.

The top features selected for the luxury score calculation include:

- `luxury_brand_indicator.0`
- `weight_to_size_ratio`
- `size_index`
- `power_to_weight_ratio`
- `luxury_brand_indicator.1`
- `compression_efficiency`
- `performance_index`
- `fuel.type.diesel`
- `bore`
- `fuel.system.mphi`

This set of features effectively captures various elements of performance, branding, size, and engine characteristics that contribute to the overall luxury perception of a vehicle.

## Clustering Results

### Cluster 1: Mid-Range Practical Cars

The luxury score for this cluster is moderate, with a mean of 0.43, indicating a blend of basic and mid-range features. The price range is moderate to high, with a mean of 1.06, reflecting mid-range pricing.

- **Performance and Power:** The Power-to-Weight Ratio is high with a mean of 1.08, offering decent acceleration and handling, while the Performance Index is moderate at a mean of 0.41, indicating a reasonable focus on performance.
- **Design and Practicality:** The Size Index is moderate with a mean of 0.37, indicating balanced car sizes suitable for daily use. The Weight-to-Size Ratio is high at 1.12, reflecting heavier builds relative to their size, and the Doors-to-Weight Ratio is low at -0.66, suggesting a preference for practicality over design complexity.
- **Luxury Indicator:** The Luxury Brand Indicator.1 has a high mean of 0.97, reflecting a strong presence of mid-range vehicles from luxury brands.

- **Fuel Economy and Efficiency:** The Fuel Economy Difference is slightly below average with a mean of -0.11, indicating adequate fuel efficiency, and Compression Efficiency is near average at -0.03, showcasing typical engine tuning.

Cluster 1 represents practical, mid-range vehicles that focus on offering a balanced driving experience. These cars are designed for buyers seeking reliable, moderately priced vehicles with good performance and practicality. While some luxury brands are represented, the overall focus is not on luxury but rather on value and usability.

## Cluster 2: High-Performance Large Cars

The luxury score for this cluster is high, with a mean of 0.84, indicating strong luxury features. The price range is moderate to high, with a mean of 0.63, reflecting higher-end pricing due to size and performance.

- **Performance-Focused:** The Power-to-Weight Ratio is low with a mean of -1.09, emphasizing fuel efficiency over raw power. Compression Efficiency is very high at 2.73, showcasing advanced engine tuning, and the Performance Index is also very high at 2.44, indicating a strong focus on high performance.
- **Design and Size:** The Size Index is large with a mean of 1.35, reflecting significant car sizes, while the Doors-to-Weight Ratio is balanced at 0.04, suggesting practical design.
- **Fuel Economy:** The Fuel Economy Difference is very low with a mean of -1.27, indicating poor fuel efficiency.
- **Luxury Indicator:** The Luxury Brand Indicator.1 has a high mean of 0.41, reflecting a strong presence of luxury-branded cars.

Cluster 2 vehicles are designed for performance-oriented buyers who value size, cutting-edge technology, and luxury. These vehicles likely represent premium SUVs or large sports sedans.

## Cluster 3: Economical Compact Cars

The luxury score for this cluster is low, with a mean of -0.27, indicating limited luxury features. The price range is low, with a mean of -0.52, making these the most economical cars.

- **Compact and Lightweight:** The Size Index is small with a mean of -0.32, reflecting compact car sizes, and the Weight-to-Size Ratio is low at -0.51, suggesting lightweight designs. The Doors-to-Weight Ratio is high with a mean of 0.27, indicating practical designs.
- **Fuel Economy and Efficiency:** The Fuel Economy Difference is above average with a mean of 0.20, reflecting better fuel efficiency, while Compression Efficiency is below average at -0.32, indicating less advanced engine tuning.
- **Performance:** The Power-to-Weight Ratio is below average with a mean of -0.33, reflecting moderate performance, while the Performance Index is low at -0.47, showing limited focus on performance.
- **Luxury Indicator:** The Luxury Brand Indicator.1 has a low mean of -0.46, reflecting minimal luxury branding.

Cluster 3 cars appeal to cost-conscious buyers seeking compact, practical, and fuel-efficient vehicles with limited luxury features. These cars represent affordable, small-sized offerings with moderate performance.



## Summary of Clusters

- **Cluster 1:** Mid-range practical cars offering reliability and balanced features at moderate pricing.
- **Cluster 2:** High-performance, large luxury cars with advanced engine technology and moderate fuel efficiency.
- **Cluster 3:** Economical compact cars with practical designs, better fuel efficiency, and limited luxury branding.

The distribution of luxury scores across the clusters reveals significant differentiation in perceived luxury, as illustrated in Figure 7. The radar chart in Figure 8 visually compares the clusters across key features, highlighting their unique characteristics and relative strengths. Figure 9 compares the mean values of selected features across the three clusters, providing insight into the unique characteristics of each group.

## Business Insights

The luxury score and selected features provide car brands with insights to refine their products. By understanding which features impact luxury perception the most, brands can strategically enhance those areas. This can be particularly useful for product development or repositioning a model to better meet market demands.

Clustering analysis helps car brands identify distinct market segments, enabling tailored marketing and product design strategies. For instance, brands can focus on fuel efficiency for economical clusters or enhance luxury features for high-end segments.

For Toyota, these insights are directly applicable. By benchmarking their high-end models, such as the Lexus LS, against luxury competitors like BMW and Mercedes, Toyota can identify gaps in features like size index and weight-to-size ratio. Improving these aspects can boost the luxury perception of Lexus vehicles. Additionally, clustering insights help Toyota understand its market position and identify opportunities to align Lexus more closely with premium segments.

## Conclusion

This project combines the development of a **Luxury Score** and **clustering analysis for market segmentation** to provide a data-driven framework for understanding and quantifying automotive luxury. The Luxury Score translates subjective perceptions of luxury into a measurable index by leveraging features most strongly associated with car pricing. It offers actionable insights for product enhancement and market positioning. Clustering analysis further segments car models into three distinct groups—mid-range practical cars, high-performance luxury vehicles, and economical compact cars—highlighting key market dynamics and customer preferences. Together, these methodologies empower automotive brands to refine strategies, optimize offerings, and effectively target diverse consumer segments.

# Appendix

Table 1: Data Dictionary for Selected Features

Feature Name	Description
power_to_weight_ratio	Ratio of a car's power output to its weight, indicating acceleration capability.
size_index	A measure representing the overall size of the car.
luxury_brand_indicator.1	Binary indicator (1 if luxury brand, 0 otherwise).
fuel_economy_difference	Difference in fuel economy compared to a baseline value.
performance_index	An index indicating the performance capability of the car.
weight_to_size_ratio	Ratio of a car's weight to its size, indicating build compactness.
compression_efficiency	Engine's compression efficiency, reflecting power and fuel efficiency balance.
doors_to_weight_ratio	Ratio of the number of doors to the car's weight, indicating design compactness.
Luxury_Score	Composite score representing the luxury characteristics of the car.
price	Price of the car in normalized units.
luxury_brand_indicator.0	Binary indicator (0 if non-luxury brand, 1 otherwise).
fuel.type.diesel	Binary indicator (1 if diesel fuel type, 0 otherwise).
bore	Diameter of the cylinder bore in the engine, affecting engine displacement.
fuel.system.mpfi	Indicator for multi-point fuel injection system.

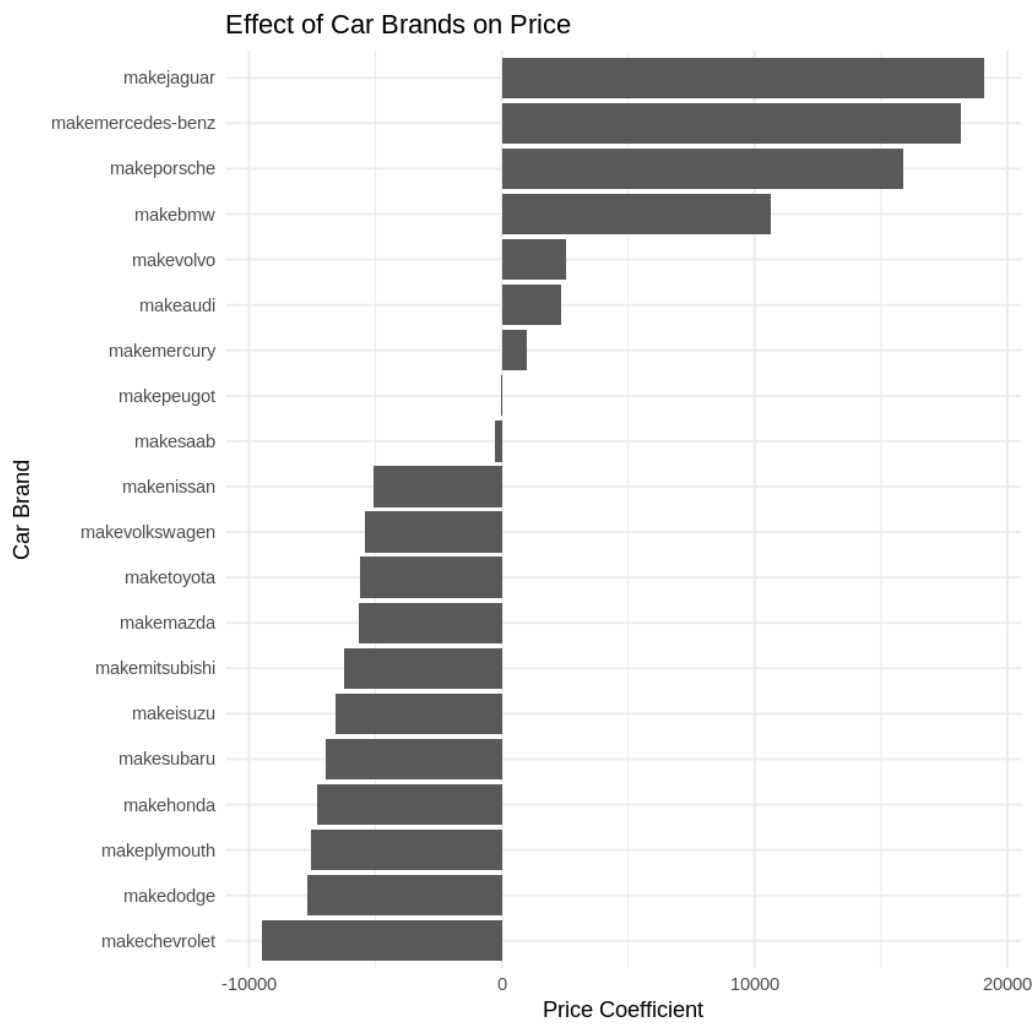


Figure 1: Effect of Car Brands on Price

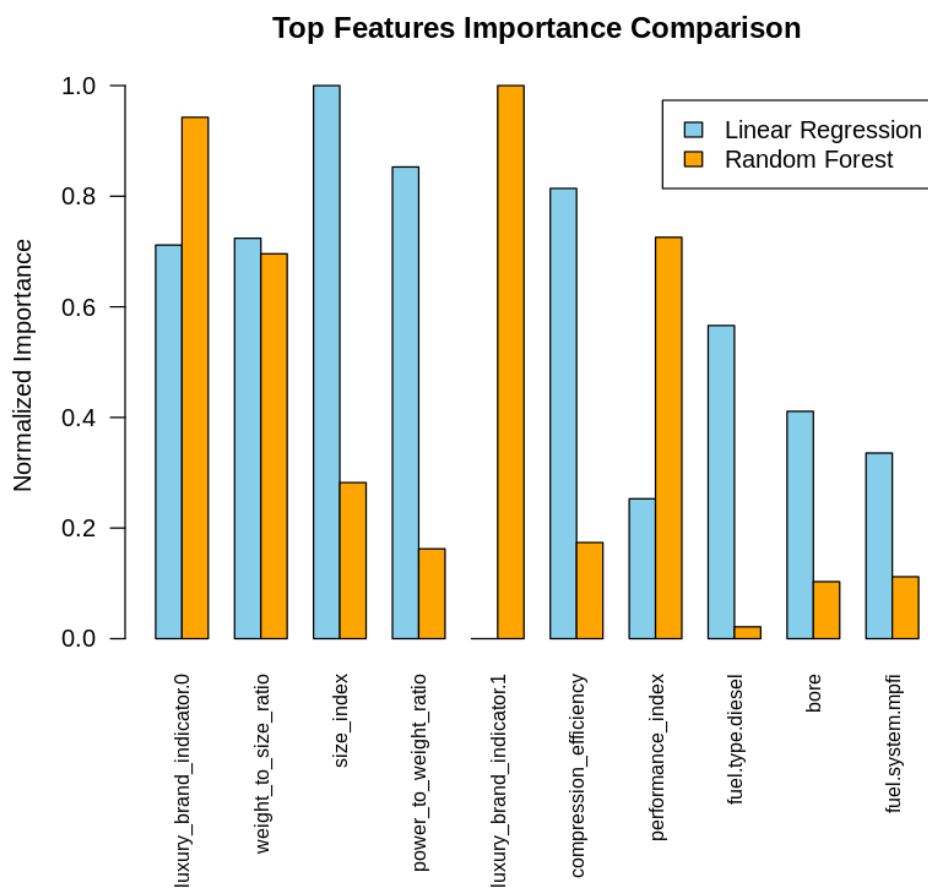


Figure 2: Top 10 Features Selected by Average Normalized Importance Scores

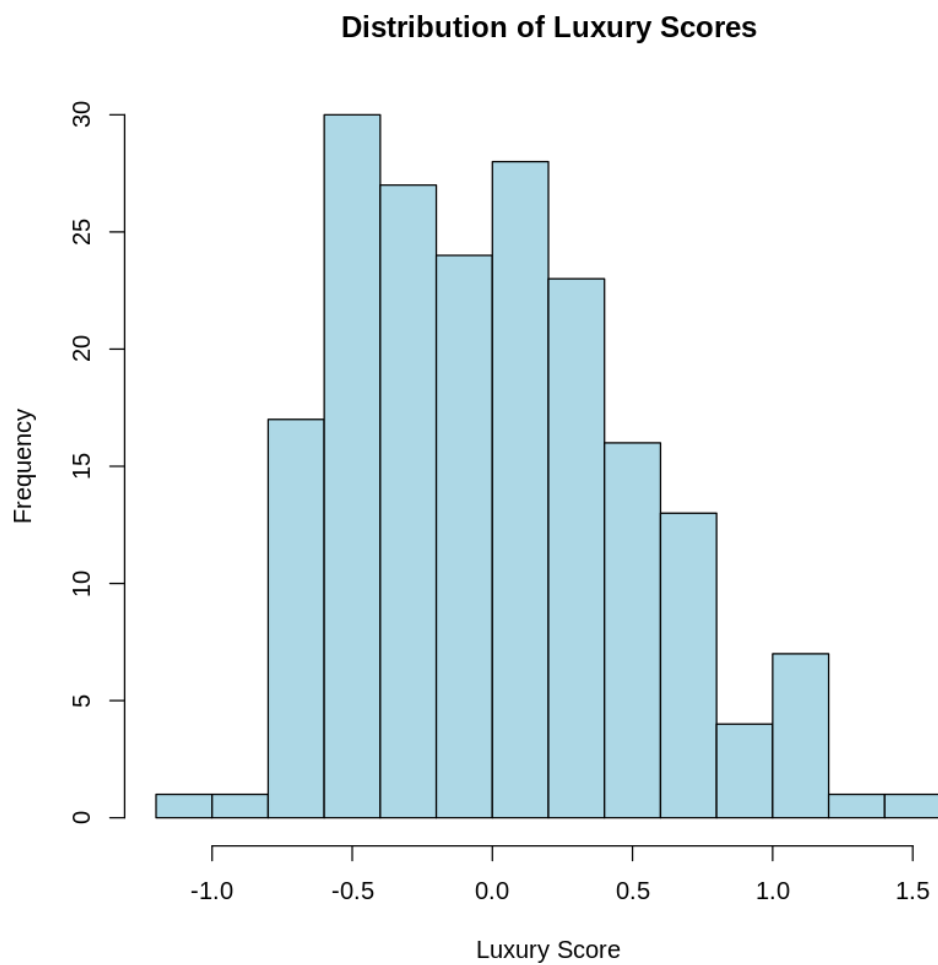


Figure 3: Distribution of Luxury Scores across all car models. This histogram visualizes how the luxury scores are distributed within the dataset, showing the clustering of scores.



Figure 4: Relationship between Luxury Score and Price. This scatter plot illustrates the correlation between a car's luxury score and its price, highlighting how luxury metrics relate to pricing.

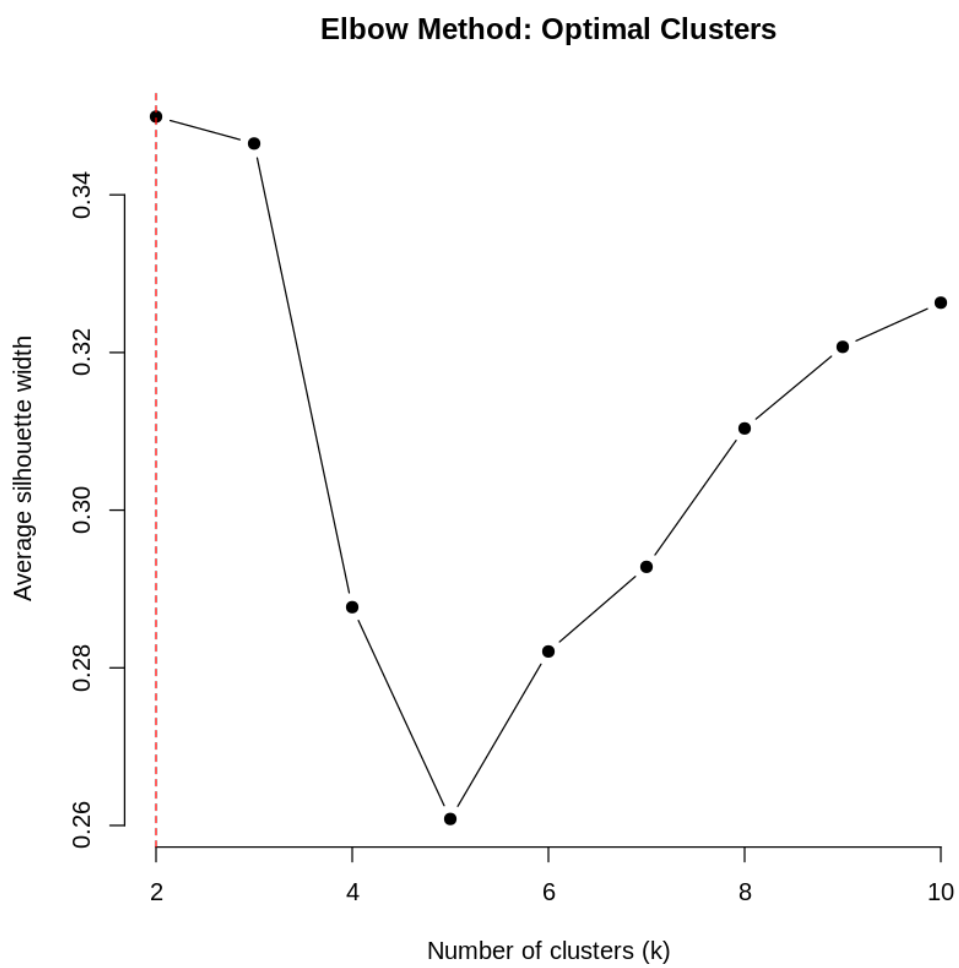


Figure 5: Silhouette Analysis for Optimal Number of Clusters

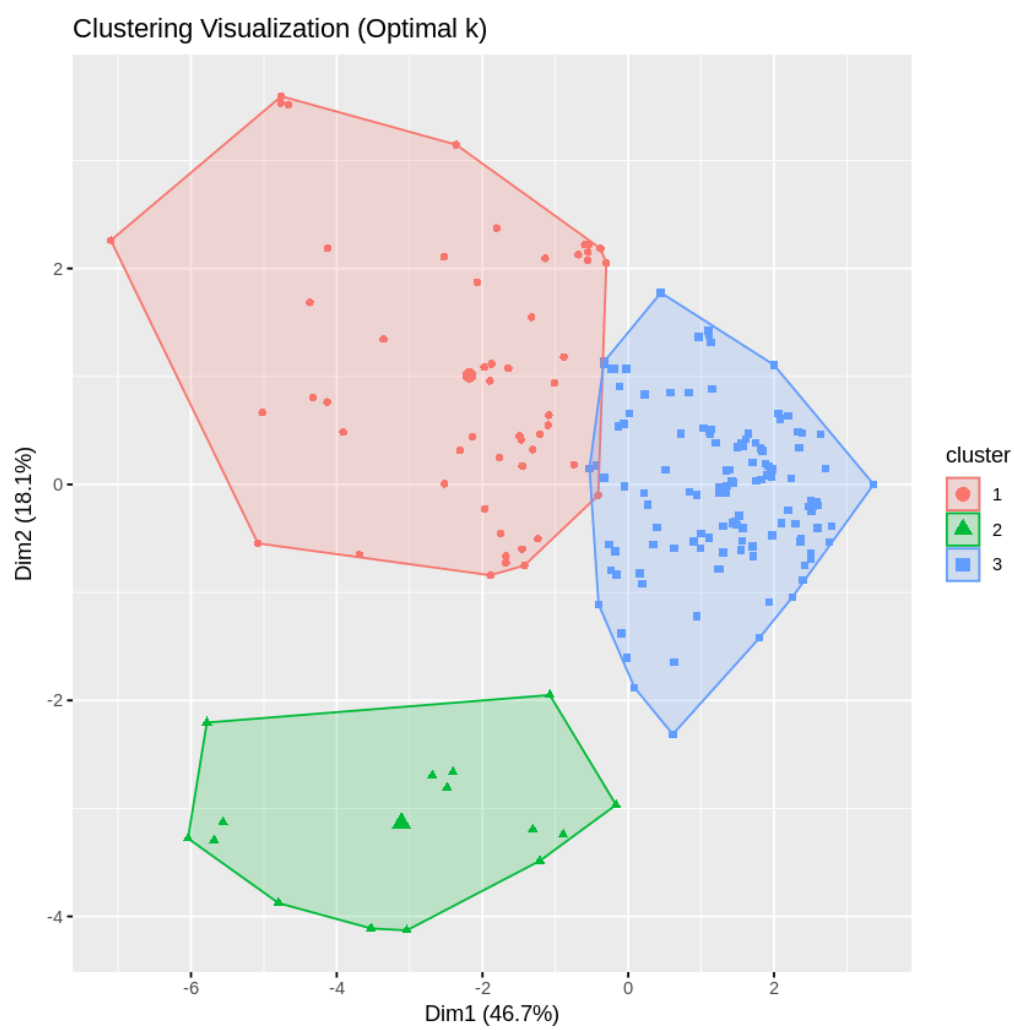


Figure 6: Visualization of Clustering Results with Optimal  $k$



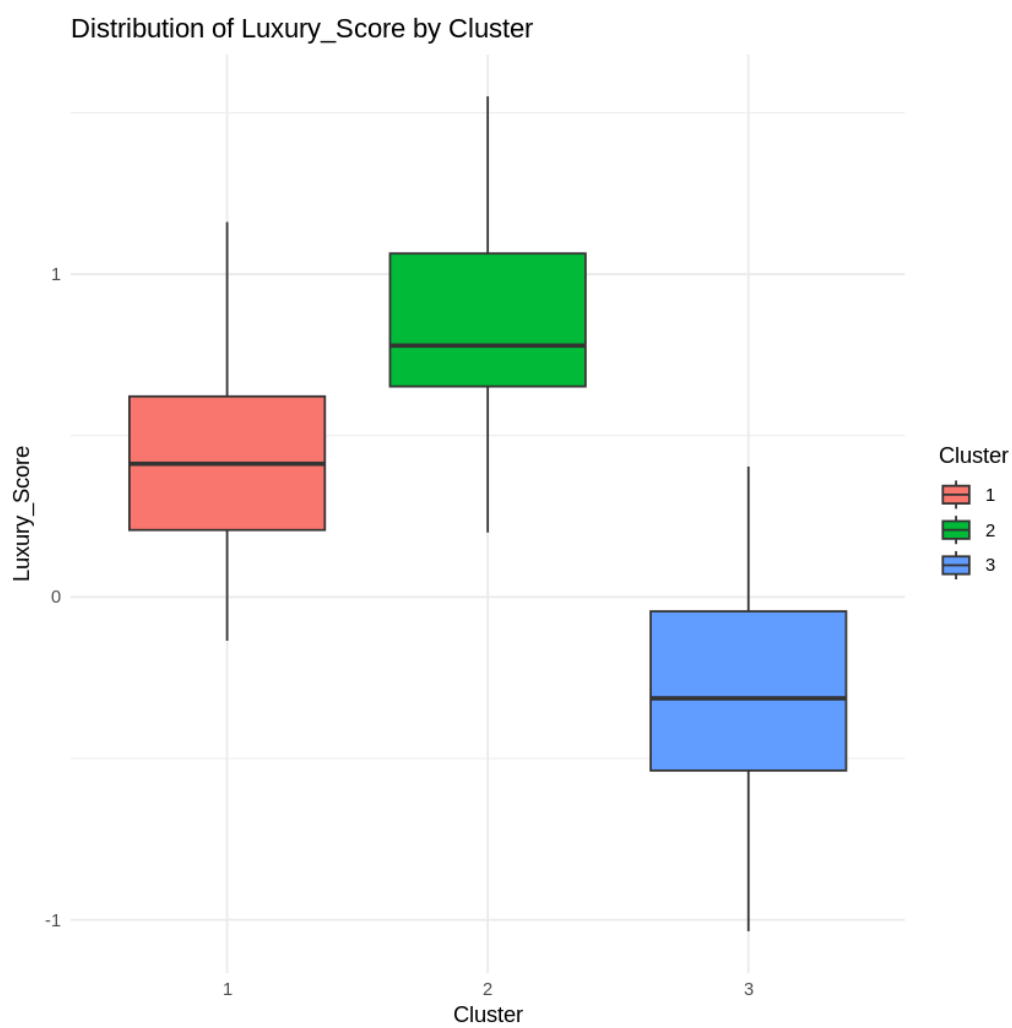


Figure 7: Distribution of Luxury Scores by Cluster. This boxplot highlights the variation in luxury scores across the three clusters, reflecting distinct levels of perceived luxury.

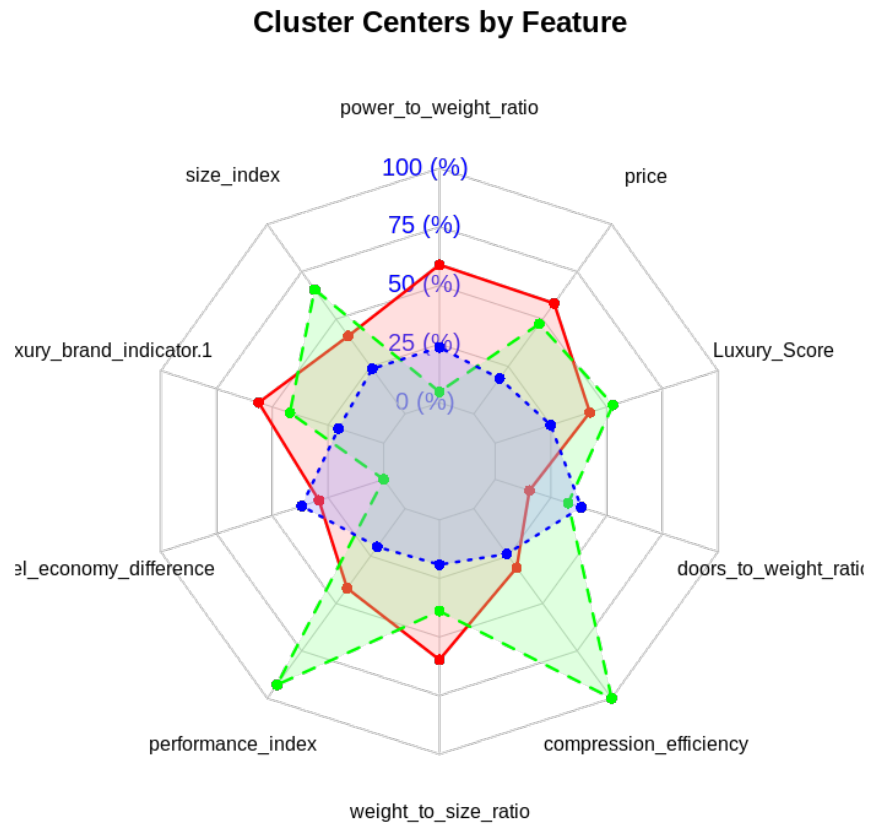


Figure 8: Radar Chart of Clusters Across Key Features. This chart illustrates the relative feature values for each cluster, providing a clear comparison of their profiles.

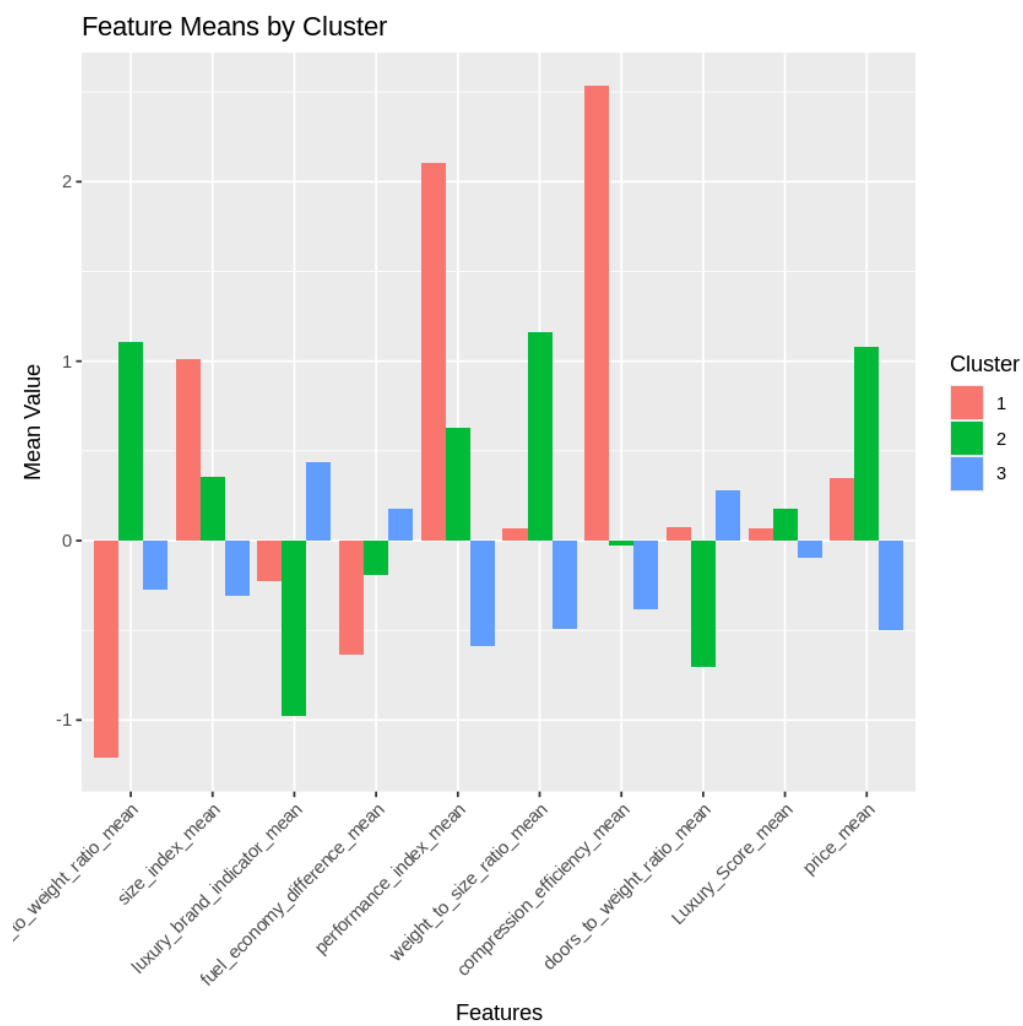


Figure 9: Feature Means by Cluster. This bar chart compares the average values of key features across clusters, highlighting their distinguishing characteristics.

## Appendix: R Code

To see more structured code with output, check my Jupyter Notebook at this link:  
<https://colab.research.google.com/drive/1brAXXqhMLEVESfciPNT3NgpVmBgJheZw?usp=sharing>.

```
1 # -*- coding: utf-8 -*-
2 # ""Automobile.ipynb
3
4 # -----
5 # Load Dataset and Packages
6 # -----
7 install.packages("dplyr")
8 install.packages("car")
9 install.packages("e1071")
10 install.packages("caret")
11 install.packages("cluster")
12 install.packages("factoextra")
13 install.packages("fmsb")
14
15 df = read.csv('Dataset5_Automobile_data.csv')
16
17 library(fmsb)
18 library(dplyr)
19 library(car)
20 library(e1071)
21 library(caret)
22 library(cluster)
23 library(factoextra)
24
25 attach(df)
26
27 # -----
28 # Data Preprocessing
29 # -----
30
31 ## Missing Values Handling
32 missing_counts <- sapply(df, function(col) sum(col == "?" | is.na(col)
33   | col == "", na.rm = TRUE))
34 missing_df <- data.frame(Column = names(missing_counts), Missing_Count
35   = missing_counts)
36 missing_df
37
38 # Replace "?" with NA for easier handling
39 df[df == "?"] <- NA
40
41 # Remove rows with any missing values
42 df <- na.omit(df)
43
44 # -----
45 # Drop Irrelevant Features
46 # -----
47 df <- df[, !(names(df) %in% c("normalized.losses", "symboling"))]
```

```

47 # -----
48 # Data Type Alignment
49 # -----
50 str(df)
51
52 ## Mapping for num.of.doors and num.of.cylinders
53 door_mapping <- c("two" = 2, "four" = 4)
54 df$num.of.doors <- door_mapping[df$num.of.doors]
55
56 cylinder_mapping <- c("two" = 2, "three" = 3, "four" = 4, "five" = 5,
57   "six" = 6, "eight" = 8, "twelve" = 12)
58 df$num.of.cylinders <- cylinder_mapping[df$num.of.cylinders]
59
60 ## Convert Columns from Character to Numeric
61 cols_to_convert <- c("bore", "stroke", "horsepower", "peak.rpm", "
  price")
62 df[cols_to_convert] <- lapply(df[cols_to_convert], function(x) as.
  numeric(as.character(x)))
63
64 # -----
65 # Feature Engineering
66 # -----
67
68 ## Power-to-Weight Ratio
69 df <- df %>%
70   mutate(power_to_weight_ratio = horsepower / 'curb.weight')
71
72 ## Size Index
73 df <- df %>%
74   mutate(size_index = 'wheel.base' * width * height)
75
76 ## Luxury Brand Indicator
77 df$make <- as.factor(df$make)
78 model <- lm(price ~ make, data = df)
79 summary(model)
80
81 luxury_brands <- c("jaguar", "mercedes-benz", "porsche", "bmw", "volvo
  ")
82 df <- df %>%
83   mutate(luxury_brand_indicator = ifelse(tolower(make) %in% luxury_
84     brands, 1, 0)) %>%
85   mutate(luxury_brand_indicator = as.factor(luxury_brand_indicator))
86
87 ## Fuel Economy Difference
88 df <- df %>%
89   mutate(fuel_economy_difference = highway.mpg - city.mpg)
90
91 ## Performance Index
92 df <- df %>%
93   mutate(performance_index = (engine.size * 'compression.ratio' * peak
94     .rpm) / 1000)
95
96 ## Weight-to-Size Ratio

```

```

94 df <- df %>%
95   mutate(weight_to_size_ratio = 'curb.weight' / size_index)
96
97 ## Compression Efficiency
98 df <- df %>%
99   mutate(compression_efficiency = ('compression.ratio' * horsepower) /
100     engine.size)
101
102 ## Doors-to-Weight Ratio
103 df <- df %>%
104   mutate(doors_to_weight_ratio = 'num.of.doors' / 'curb.weight')
105
106 # -----
107 # Multicollinearity & Feature Selection
108 # -----
109
110 ## Drop Redundant Features
111 df <- df %>%
112   select(-wheel.base, -length, -width, -curb.weight, -engine.size, -
113     highway.mpg, -city.mpg, -horsepower, -compression.ratio, -make)
114
115 ## Calculate VIF Scores
116 numeric_vars <- sapply(df, is.numeric)
117 numeric_df <- df[, numeric_vars]
118 vif_scores <- vif(lm(price ~ ., data = numeric_df))
119 vif_df <- data.frame(Variable = names(vif_scores), VIF = vif_scores)
120 vif_df <- vif_df[order(-vif_df$VIF), ]
121 vif_df
122
123 # -----
124 # Outlier Analysis
125 # -----
126
127 ## Detect Outliers Using IQR Method
128 detect_outliers <- function(x) {
129   q1 <- quantile(x, 0.25)
130   q3 <- quantile(x, 0.75)
131   iqr <- q3 - q1
132   lower_bound <- q1 - 1.5 * iqr
133   upper_bound <- q3 + 1.5 * iqr
134   outliers <- x[x < lower_bound | x > upper_bound]
135   return(length(outliers))
136 }
137
138 outlier_counts <- sapply(df[, sapply(df, is.numeric)], detect_outliers)
139
140 outlier_df <- data.frame(Column = names(outlier_counts), Outlier_Count
141   = outlier_counts)
142 outlier_df
143
144 # -----
145 # Boxplots for Selected Features
146 # -----

```

```

143 par(mfrow = c(1, 2))
144 boxplot(df$fuel_economy_difference, main = "Fuel_Economy_Difference",
145         ylab = "Fuel_Economy_Difference")
146 boxplot(df$performance_index, main = "Performance_Index", ylab = "
147         Performance_Index")
148
149 # -----
150 # Dummification of Character/Factor Columns
151 # -----
152
153 ## Dummify All Character/Factor Variables
154 char_factor_cols <- names(df)[sapply(df, function(x) is.character(x) |
155                                     is.factor(x))]
156 for (col in char_factor_cols) {
157   if (is.character(df[[col]])) {
158     df[[col]] <- as.factor(df[[col]])
159   }
160   dummy_vars <- dummyVars(paste("~", col), data = df)
161   dummy_df <- data.frame(predict(dummy_vars, newdata = df))
162   df <- cbind(df, dummy_df)
163   df <- df[, -which(names(df) == col)]
164 }
165
166 # -----
167 # Standardization
168 # -----
169
170 numeric_cols <- sapply(df, is.numeric)
171 df_numeric <- df[, numeric_cols]
172
173 # Standardize the numeric columns
174 df_scaled <- scale(df_numeric)
175
176 # Convert scaled data back to data frame
177 df_scaled <- as.data.frame(df_scaled)
178
179 str(df_scaled)
180
181 # -----
182 # Feature Selection
183 # -----
184
185 ## Linear Regression Model
186 lm_model <- lm(price ~ ., data = df_scaled)
187
188 # Extract coefficients as feature importance
189 lm_importance <- summary(lm_model)$coefficients[, "Estimate"]
190 lm_importance <- data.frame(Feature = names(lm_importance), Importance
191                             = lm_importance)
192 print(lm_importance)
193
194 ## Random Forest Model
195
196 # Load Random Forest Library

```

```

192 library(randomForest)
193
194 # Fit Random Forest Model
195 rf_model <- randomForest(price ~ ., data = df_scaled, importance =
  TRUE)
196
197 # Extract Feature Importance
198 rf_importance <- data.frame(Feature = rownames(rf_model$importance),
199                             Importance = rf_model$importance[, "
  IncNodePurity"])
200 print(rf_importance)
201
202 # Plot Feature Importance
203 barplot(rf_importance$Importance, names.arg = rf_importance$Feature,
  las = 2, main = "Random_Forest_Feature_Importance")
204
205 # -----
206 # Final Feature Selection and Normalization
207 # -----
208
209 # Normalize a Vector to Range 0-1
210 normalize <- function(x) {
211   return((x - min(x)) / (max(x) - min(x)))
212 }
213
214 # Normalize Feature Importance from Linear Regression and Random
  Forest
215 lm_importance$Normalized_Importance <- normalize(abs(lm_importance$
  Importance))
216 rf_importance$Normalized_Importance <- normalize(rf_importance$
  Importance)
217
218 # Rename Columns for Clarity Before Merging
219 colnames(lm_importance) <- c("Feature", "Linear_Importance", "Linear_
  Normalized")
220 colnames(rf_importance) <- c("Feature", "RandomForest_Importance", "
  RandomForest_Normalized")
221
222 # Merge Importance from Both Models by Feature
223 combined_importance <- merge(lm_importance, rf_importance, by = "
  Feature", all = TRUE)
224
225 # Fill NA Values with 0 (In Case a Feature is Missing from One of the
  Models)
226 combined_importance[is.na(combined_importance)] <- 0
227
228 # Calculate Average Normalized Importance Across Methods
229 combined_importance$Average_Normalized_Importance <- rowMeans(
230   combined_importance[, c("Linear_Normalized", "RandomForest_
    Normalized")])
231
232 # Sort by Average Normalized Importance
233 combined_importance <- combined_importance[order(-combined_importance$

```



```

    Average_Normalized_Importance), ]
234
235 # Select Top N Features
236 top_features <- head(combined_importance, 10)
237
238 # Print Combined Importance Table and Top Features
239 print(combined_importance)
240 print(top_features)
241
242 # Adjust Plot Size and Margins for Visualization
243 par(mar = c(12, 5, 4, 2))
244
245 # Plot Top Features
246 barplot(
247   height = t(as.matrix(top_features[, c("Linear_Normalized", "
248     RandomForest_Normalized")])),
249   beside = TRUE,
250   names.arg = top_features$Feature,
251   las = 2,
252   col = c("skyblue", "orange"),
253   legend.text = c("Linear_Regression", "Random_Forest"),
254   main = "Top_Features_Importance_Comparison",
255   ylab = "Normalized_Importance",
256   cex.names = 0.8
257 )
258
259 # -----
260 # Luxury Score Construction
261 # -----
262
263 ## Luxury Score Calculation
264 calculate_luxury_score <- function(df, top_features) {
265   # Extract Top Features and Their Normalized Importance
266   selected_features <- top_features$Feature
267   feature_weights <- top_features$Average_Normalized_Importance
268
269   # Ensure Feature Weights Sum to 1
270   feature_weights <- feature_weights / sum(feature_weights)
271
272   # Create Luxury Score Column
273   df$Luxury_Score <- rowSums(df[, selected_features] * feature_weights
274     )
275
276   return(df)
277 }
278
279 # Apply Function to Scaled Dataset
280 df_scaled_with_luxury_score <- calculate_luxury_score(df_scaled, top_
281   features)
282
283 # Plot Distribution of Luxury Scores
284 hist(
285   df_scaled_with_luxury_score$Luxury_Score,

```

```

283   main = "Distribution of Luxury Scores",
284   xlab = "Luxury Score",
285   col = "lightblue",
286   border = "black",
287   breaks = 15
288 )
289
290 # Scatter Plot: Luxury Score vs Price
291 plot(
292   df_scaled_with_luxury_score$Luxury_Score,
293   df_scaled_with_luxury_score$price,
294   main = "Luxury Score vs Price",
295   xlab = "Luxury Score",
296   ylab = "Price",
297   col = "darkblue",
298   pch = 19
299 )
300
301 # -----
302 # Clustering Analysis
303 # -----
304
305 ## Subset Dataset for Clustering
306 clustering_features <- df_scaled_with_luxury_score[, c(
307   "power_to_weight_ratio",
308   "size_index",
309   "luxury_brand_indicator.1",
310   "fuel_economy_difference",
311   "performance_index",
312   "weight_to_size_ratio",
313   "compression_efficiency",
314   "doors_to_weight_ratio",
315   "Luxury_Score",
316   "price"
317 )]
318
319 # -----
320 # Silhouette Analysis to Determine Optimal Clusters
321 # -----
322
323 silhouette_analysis <- function(data, max_clusters = 10) {
324   sil_width <- numeric(max_clusters - 1)
325
326   for (k in 2:max_clusters) {
327     kmeans_model <- kmeans(data, centers = k, nstart = 25)
328     sil <- silhouette(kmeans_model$cluster, dist(data))
329     sil_width[k - 1] <- mean(sil[, 3])
330   }
331
332   return(sil_width)
333 }
334
335 # Perform Silhouette Analysis

```

```

336 max_clusters <- 10
337 sil_width <- silhouette_analysis(clustering_features, max_clusters)
338
339 # Plot Silhouette Scores for Each k
340 plot(2:max_clusters, sil_width, type = "b", pch = 19, frame = FALSE,
341      xlab = "Number_of_clusters_(k)", ylab = "Average_silhouette_width",
342      main = "Elbow_Method:_Optimal_Clusters")
343 abline(v = which.max(sil_width) + 1, col = "red", lty = 2)
344
345 # Determine Optimal Number of Clusters
346 optimal_k <- which.max(sil_width) + 1
347 cat("Optimal_number_of_clusters_based_on_silhouette_score:", optimal_k
348     , "\n")
349 # -----
350 # K-Means Clustering with Optimal Number of Clusters
351 # -----
352
353 final_kmeans <- kmeans(clustering_features, centers = 3, nstart = 25)
354
355 # Visualize Clustering Results
356 library(factoextra)
357 fviz_cluster(final_kmeans, data = clustering_features, geom = "point",
358              main = "Clustering_Visualization_(Optimal_k)")
359
360 # -----
361 # Cluster Analysis
362 # -----
363
364 # Add Cluster Labels to Dataset
365 df_scaled_with_luxury_score$cluster <- final_kmeans$cluster
366
367 # Calculate Cluster-Level Statistics
368 cluster_stats <- df_scaled_with_luxury_score %>%
369   group_by(cluster) %>%
370   summarize(
371     mean_luxury_score = mean(Luxury_Score),
372     median_luxury_score = median(Luxury_Score),
373     sd_luxury_score = sd(Luxury_Score),
374     count = n()
375   )
376
377 print(cluster_stats)
378
379 # -----
380 # Visualization: Boxplot of Feature Distribution by Cluster
381 # -----
382
383 # Select Feature for Boxplot (e.g., 'Luxury_Score')
384 feature_to_plot <- "Luxury_Score"
385
386 # Create Boxplot for Selected Feature by Cluster

```

```

387 ggplot(df_scaled_with_luxury_score, aes(x = factor(cluster), y = .data
    [[feature_to_plot]], fill = factor(cluster))) +
388   geom_boxplot() +
389   labs(title = paste("Distribution of", feature_to_plot, "by Cluster")
    ,
390     x = "Cluster",
391     y = feature_to_plot,
392     fill = "Cluster") +
393   theme_minimal()
394
395 # -----
396 # Radar Chart of Cluster Centers by Feature
397 # -----
398
399 cluster_centers <- as.data.frame(final_kmeans$centers)
400 cluster_centers <- rbind(rep(max(cluster_centers), ncol(cluster_
    centers)),
401                           rep(min(cluster_centers), ncol(cluster_
    centers)),
402                           cluster_centers)
403
404 radarchart(cluster_centers,
405             axistype = 1,
406             pcol = c("red", "green", "blue"),
407             pfcol = adjustcolor(c("#FF9999", "#99FF99", "#9999FF"),
    alpha.f = 0.3),
408             plwd = 2,
409             cglcol = "grey", cglty = 1, cglwd = 0.8,
410             vlce = 0.8,
411             title = "Cluster Centers by Feature")
412
413 # -----
414 # Cluster-Level Statistics for Specified Variables
415 # -----
416
417 cluster_stats <- df_scaled_with_luxury_score %>%
418   group_by(cluster) %>%
419   summarize(
420     across(
421       c(
422         "power_to_weight_ratio",
423         "size_index",
424         "luxury_brand_indicator.1",
425         "fuel_economy_difference",
426         "performance_index",
427         "weight_to_size_ratio",
428         "compression_efficiency",
429         "doors_to_weight_ratio",
430         "Luxury_Score",
431         "price"
432       ),
433       list(mean = mean)
434     ),

```

```

435     count = n()
436   )
437
438 print(cluster_stats)
439
440 # -----
441 # Visualization: Feature Means by Cluster
442 # -----
443
444 cluster_means <- tibble::tibble(
445   cluster = c(1, 2, 3),
446   power_to_weight_ratio_mean = c(-1.2112051, 1.1070278, -0.2719148),
447   size_index_mean = c(1.0100569, 0.3585888, -0.3047082),
448   luxury_brand_indicator_mean = c(-0.2277566, -0.9761798, 0.4399394),
449   fuel_economy_difference_mean = c(-0.6336319, -0.1890356, 0.1762587),
450   performance_index_mean = c(2.1013806, 0.6273455, -0.5847224),
451   weight_to_size_ratio_mean = c(0.06967552, 1.16250959, -0.49277906),
452   compression_efficiency_mean = c(2.53229070, -0.02752784,
453     -0.37975287),
454   doors_to_weight_ratio_mean = c(0.07208247, -0.70114069, 0.27958218),
455   Luxury_Score_mean = c(0.06761755, 0.18050763, -0.09856919),
456   price_mean = c(0.3484391, 1.0797236, -0.5015142)
457 )
458
459 cluster_means_melted <- melt(cluster_means, id.vars = "cluster")
460
461 ggplot(cluster_means_melted, aes(x = variable, y = value, fill =
462   factor(cluster))) +
463   geom_bar(stat = "identity", position = "dodge") +
464   labs(title = "Feature Means by Cluster",
465     x = "Features",
466     y = "Mean Value",
467     fill = "Cluster") +
468   theme(axis.text.x = element_text(angle = 45, hjust = 1))

```

Listing 1: Quantify Luxury and Clustering for Market Segmentation