# Indian  Automobile Strategic Grouping

**CSE 3506 Essentials of Data Analytics J-component**

**Team Members**

Aakash R 20BCE1003

Subramanian Nachiappan 20BCE1019

Komma Reddy Saketh Reddy 20BCE1042

Sanjeet V 20BCE1050

**Under the Guidance of**

Dr. Asnath Victy Phamila Y

# Table of Content

# Chapter 1 - Introduction

## 1.1 Abstract

The "Indian Automotive Strategic Grouping" is a study that aims to understand and analyse the Indian automotive market. The study utilises a data-driven approach, grouping different automobile models into distinct strategic groups based on various attributes such as price, brand, segment, and features. Here we use a clustering method, to identify different strategic groups of car models in the automobile industry. Clustering or cluster analysis is an unsupervised machine learning technique, which groups the unlabelled dataset. The grouping of vehicles into strategic clusters helps to understand the market dynamics and provides insights into consumer preferences and buying behaviour.The findings of this study can be useful for automobile manufacturers, dealers, and marketers to develop targeted marketing strategies, improve product offerings, and better understand consumer demand.The study also provides valuable insights into the Indian automotive market and its trends, and can be used as a reference for future research and analysis in this field. The aim of the project is also to find out the best car model in each cluster based on value for money and the features it has and also suggest users a cluster that most matches his/her choice of attributes.

.

## 1.2 Keywords:

Automobile, Segmentation, Strategic Grouping, Competitors, Unsupervised Learning, Clustering

## 1.3 Introduction

The "Indian Automotive Strategic Grouping" study is a comprehensive analysis of the Indian automotive market that utilises a data-driven approach to understand consumer preferences and buying behaviour. This study utilises a clustering method, a popular unsupervised machine learning technique, to group automobile models into distinct strategic groups based on various attributes, including price, brand, segment, and features. By identifying these strategic groups, the study provides valuable insights for automobile manufacturers, dealers, and marketers to develop targeted marketing strategies and improve product offerings. Additionally, the study identifies the best car model in each cluster based on value for money and features, and suggests to users a cluster that best matches their choice of attributes. The findings of this study can be used as a reference for future research and analysis in the field of Indian automotive market trends. This study is a must-read for industry professionals and anyone interested in gaining a deeper understanding of the Indian automotive market.

The Indian automotive industry comprises various segments, including passenger vehicles, commercial vehicles, two-wheelers, three-wheelers, and tractors. The industry is highly diversified, with a mix of domestic and international players operating across different segments. Some of the major players in the **Indian automotive industry include Maruti Suzuki, Tata Motors, Mahindra & Mahindra, Hyundai, Honda, and Toyota**.

The Indian automotive industry has grown rapidly over the past few decades, with the expansion of the domestic market and the increasing adoption of new technologies. In recent years, the industry has faced several challenges, including changing consumer preferences, tightening emission norms, and the increasing competition from international players.

**Strategic Grouping in the Indian Automotive Industry:**

Strategic grouping is a concept that refers to the clustering of firms with similar market positions and competitive strategies. In the Indian automotive industry, firms can be grouped based on various criteria, such as the market segments they operate in, their product offerings, and their geographical locations. Strategic grouping has important implications for industry structure, competitive dynamics, and firm performance.

The Indian Automotive Strategic Grouping dataset provides a detailed analysis of the evolution of strategic groups in the Indian automotive industry over time. The dataset identifies the major strategic groups in the industry and analyzes the factors that determine the boundaries of these groups. The dataset also examines the competitive dynamics between firms within each group and the implications for industry structure and firm performance.

# Chapter 2: Literature Survey

**[1]** **"Vehicle price prediction system using machine learning techniques"** by K. Noor and S. Jan presented a study on using machine learning algorithms to predict vehicle prices. They compared several algorithms, including linear regression, decision trees, and support vector regression, and found that support vector regression performed the best with high accuracy and low mean absolute error. The authors also found that the make, model, and year of a vehicle were the most significant factors impacting its price, followed by engine size and other attributes. This study demonstrates the potential of machine learning techniques in predicting vehicle prices and adds valuable insights to the literature on analytics in the automobile industry.

**[2] "How much is my car worth? A methodology for predicting used car prices using random forest"** by N. Pal, P. Arora, P. Kohli, D. Sundararaman, and S. S. Palakurthy proposes a method for predicting used car prices using random forest algorithms. The authors collected data on used car attributes such as make, model, year, and condition and used this data to train the random forest model. The results of their study showed that the random forest algorithm outperformed other machine learning algorithms in terms of accuracy and mean absolute error in predicting used car prices. The authors concluded that the proposed methodology provides a reliable and efficient way to predict used car prices and can be used as a reference for industry professionals. This study adds to the growing body of literature on analytics in the automobile industry and provides valuable insights into the use of machine learning techniques in predicting used car prices.

**[3] "Different Clustering Algorithms for Big Data Analytics: A Review"** by Meenu Dave and Hemant Gianey covers a range of clustering techniques including k-means, hierarchical clustering, and density-based clustering. The study highlights the advantages and disadvantages of each algorithm and provides a comparison of their performance in different applications. The paper contributes to the existing literature on clustering algorithms for big data analytics and provides insights for researchers and practitioners seeking to select the most appropriate algorithm for their data.

**[4] "The Strategic Grouping of Automobile Firms: A Review of Literature"** by Michael E. Porter: This paper provides a comprehensive overview of strategic grouping in the automobile industry. It identifies the factors that determine the boundaries of strategic groups and how firms within these groups compete with each other. The paper also discusses the implications of strategic grouping for industry structure and competitive dynamics.

**[5] "Strategic Grouping and Competitive Dynamics in the Global Automobile Industry"** by Yu-Shan Chen and Yi-Jen Chen: This paper examines the strategic grouping of global automobile firms and their competitive dynamics. The study analyzes the evolution of strategic groups over time and how firms use various strategic actions to improve their competitive positions. The paper also discusses the impact of globalization on strategic grouping.

**[6] "Strategic Grouping in the US Automobile Industry"** by David M. Cutler and James M. Poterba: This study investigates the strategic grouping of US automobile firms based on their product offerings, market segments, and geographic locations. The paper analyzes how firms within these groups compete with each other and the implications for industry structure. The study also examines the impact of government policies on strategic grouping in the US automobile industry.

**[7] "The Effect of Strategic Grouping on Firm Performance: Evidence from the Japanese Automobile Industry"** by Masahiko Aoki and Yujiro Hayami: This paper explores the relationship between strategic grouping and firm performance in the Japanese automobile industry. The study shows that firms in strategic groups have similar performance levels and that strategic grouping can improve the efficiency of industry production processes. The paper also discusses the implications of strategic grouping for industry innovation.

[**8] "Strategic Grouping in the Chinese Automobile Industry"** by Xuezhi Qin and Jiajia Liu: This paper examines the strategic grouping of Chinese automobile firms and their competitive dynamics. The study identifies five strategic groups in the industry and analyzes their different competitive strategies and market positions. The paper also discusses the implications of strategic grouping for the future development of the Chinese automobile industry, particularly in the context of government policies and industry regulations.

**[9] "How Phases of Cluster Development are Associated with Innovation - the Case of China"** by William Fri and Tobias Pehrsson says that Clusters are geographic concentrations of interconnected businesses and institutions that promote innovation and economic growth. Studies suggest that clusters stimulate innovation through collaboration, competition, and knowledge spillovers. The different phases of cluster development, such as growth and maturity, are also associated with increased innovation. In China, clusters have been found to overcome the lack of innovation capabilities and resources in small and medium-sized enterprises, and to facilitate collaboration and access to specialised resources. The paper uses a qualitative case study approach to examine the relationship between the phases of cluster development and innovation in a high-tech cluster in Guangdong province.

**[10] "Industrial Clusters in India: Evidence from Automobile Clusters in Chennai and the National Capital Region"** by Okada Aya and Siddharthan N.S. investigates the nature and characteristics of industrial clusters in India, specifically in the automobile industry in Chennai and the National Capital Region. The study highlights the role of agglomeration economies and local institutional support, as well as the need for collaboration among firms to enhance competitiveness. The paper contributes to the existing literature on industrial clusters in India and the automobile industry, and provides insights for policymakers and businesses seeking to promote cluster development. However, limitations such as the small sample size and lack of generalizability should be considered.

**[11]** "**A latent variable model for mixed mode data: The Mixture-of-Distribution approach"** by Moustaki and Knott used Gower's distance as a distance metric in a hierarchical clustering algorithm for mixed data. The study compared the performance of the proposed algorithm with other clustering algorithms based on other distance metrics. The results showed that the proposed algorithm outperformed other clustering algorithms in terms of clustering quality.

[12] Bhattacharyya and Mukhopadhyay (2015). In paper **Cluster validation using Calinski–Barabasz index for gene expression data** published in Journal of Statistical Planning and Inference, 157, 45-55. Author used the Calinski-Barabasz index to determine the optimal number of clusters in gene expression data. The study compared the performance of the Calinski-Barabasz index with other clustering criteria, such as the gap statistic and the average silhouette width. The results showed that the Calinski-Barabasz index provided better clustering results compared to other criteria.

[13] "**Extensions to the k-means algorithm for clustering large data sets with categorical values" by** Huang, Z. (1998) proposed the k-prototype model, an extension of the k-means algorithm, to cluster mixed data that contains both numerical and categorical variables. The k-prototype model was shown to outperform other clustering algorithms in various applications such as customer segmentation and job satisfaction surveys. The model was also modified to handle missing values and outliers in the data. The k-prototype model works by calculating

the dissimilarity between objects based on their attributes, using a combination of a distance metric for numerical variables and a dissimilarity measure for categorical variables.

[14] Book by Kaufman and Rousseeuw's "**Finding groups in data: An introduction to cluster analysis"** is an influential work on cluster analysis, providing a comprehensive introduction to the subject. The book covers a wide range of clustering methods, including hierarchical clustering, partitioning algorithms, and density-based clustering. It also discusses the choice of distance measures, the assessment of clustering quality, and the interpretation of cluster results. The book has been widely cited and used as a reference in research and practical applications, such as data mining, image analysis, and market segmentation. The authors' emphasis on the importance of exploratory data analysis and understanding the underlying structure of the data has made this book a timeless resource for anyone interested in clustering.

# Chapter 3: Proposed Methodology

Our model uses clustering which is an unsupervised machine learning technique. This data-driven approach groups car models based on the attributes like brand, price, other features and identifies different strategic grouping of car models. These groups are then studied to understand consumer behaviour and market dynamics. Finally, we determine the best car for each group based on its features. We also recommend the car that matches the customer's interests by identifying the group to which customer belongs with the most matches.

## 3.1 Overall Architecture

**COLLECT**

Browse Data Repositories to collect Indian Automobile Data

**CLEAN**

Remove unwanted variables
Checking missing variables
Data imputation

**BUILD**

Feature Selection , Feature Extraction, Clustering

**RESULT**

Examine the results and use the insights to improve business

## 3.2 Modular Description

| Data Collection | Data Cleaning and Preprocessing | Exploratory Data Analysis | Feature Selection | Model Building & Insights |
|---|---|---|---|---|
| • Identify appropriate dataset | • Cleaning features<br>• Imputation<br>• Transforming Features to apt datatype | • Visualization<br>• Descriptive Statistics | • Correlation<br>• PCA<br>• Information Gain | • Gower's Distance<br>• Partitioning Around Medoids |

1. **M1: Collection and Organisation of data**

The process of gathering and analysing accurate data from various sources to find answers to research problems, trends, probabilities, etc., to evaluate possible outcomes is known as data collection. First, we find all the possible datasets that we can use for the project which may include industry reports, market surveys, government publications, and company websites. We will then analyze the data which provide detailed information on car models, their attributes, and market trends. Through discussion we identify the problem statement and finalize the questions which can be answered at end. Accordingly we will select the dataset which meets our requirements and satisfies our problem statement.

2. **M2: Data cleaning and imputation**

First we will analyse the dataset and find the attributes which can be converted to numeric by removing the units attached to it. We are removing the attributes which are having NA values greater than 10% and checking for the class of the attribute so that we can imputing it accordingly. The numerical attributes with NA values are selected and are imputed using Mean and categorical attributes with NA values are selected and imputed using the Mode. Finally we can look at the correlation matrix

for numerical attributes to get the overview of the variables and know how they are related to each other and remove the attributes which are having a high correlations.

### 3. M3: Exploratory data analytics

Exploratory Data Analysis (EDA) is an approach to analysing data using visual techniques. It is used to discover trends, patterns, or to check assumptions with the help of statistical summaries and graphical representations. It plays very important role in unsupervised learning. We need to find the various variants of the cars for clustering the data to identify its competitor. We will Explore and Visualize all the possible graphs that can be drawn.

### 4. M4: Variable selection

From the plots we have drawn in the previous module, we will infer a lot of things. For example, the length of the vehicle is directly related to the number of seats. We must decide upon which variable will be used in the construction of a clustering model, and that should in turn help in grouping the various car models accurately. We also use information gain for selecting the categorical variables that are good predictors of the car models. Principal Component Analysis is used to extract the numerical features that explain a majority of the variance in the data.

### 5. M5: Model building

In this module we construct the machine learning model using clustering, a unsupervised learning technique. We use the Silhouette technique to determine the optimal number of clusters, with gower ' s distance as the distance metric. Since the dataset contain both numerical and categorical features, we use the Partitioning Around Medoids (PAM) clustering to segment the car models.

## 3.3 Data Set

The vehicle dataset of Indian cars from Kaggle is a comprehensive collection of data regarding various car models and their specifications. The dataset includes 60 different car models and has 34 features for each car model.The dataset consists of 1267 instances and 141 columns.This dataset can be used to analyse and compare different car models based on their features and specifications. After doing appropriate preprocessing we have reduced the number of features from 141 to 34.
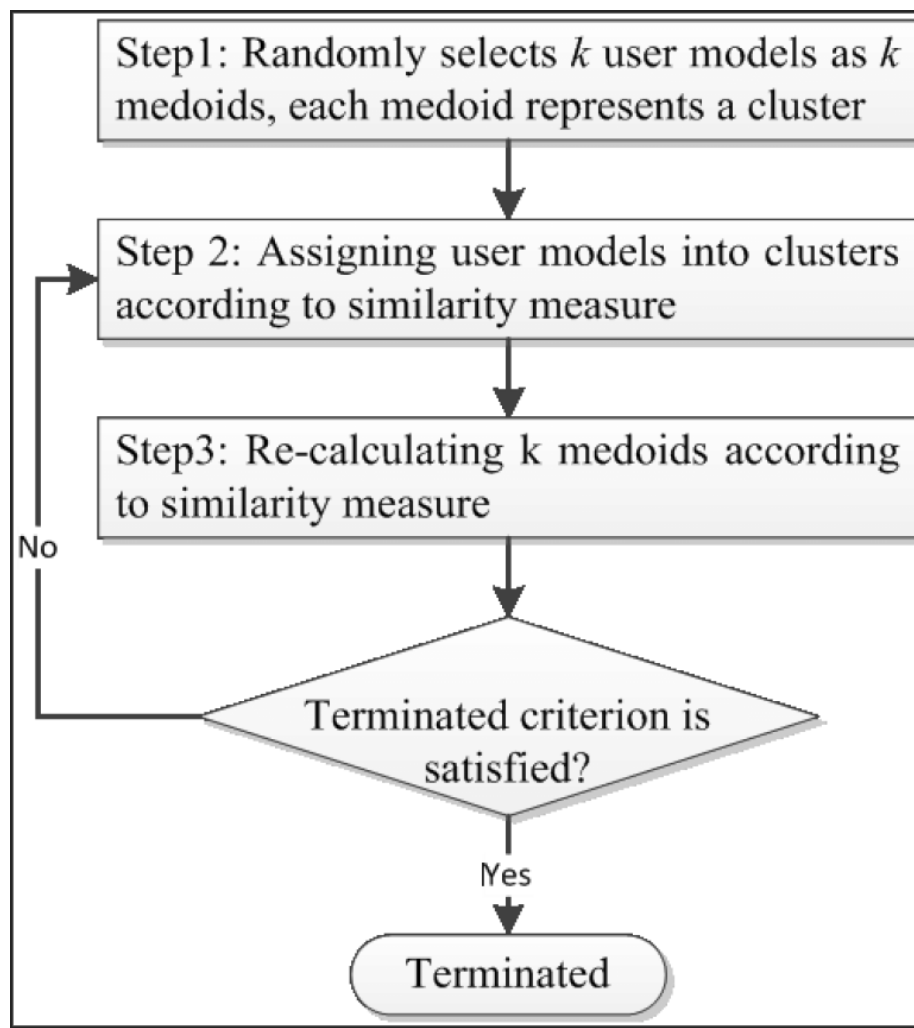
The features which had been included in the dataset are been listed below:-

- Make
- Model
- Variant
- Cylinders
- Cylinder_Config
- Fuel_Tank_Capacity
- Length
- Doors
- Rear_Brakes
- Seating_Capacity
- Instrument_Console

- Front_Brakes
- Power
- Seats_Material
- Handbrake
- Ex.Showroom_Price
- Displacement
- Valves_Per_Cylinder
- Drivetrain
- Power_Windows
- Type
- Sun_Visor

- Emission_Norm
- Engine_Location
- Fuel_Type
- Height
- Width
- Body_Type
- ARAI_Certified_Mileage
- Gears
- Torque
- Wheelbase
- Car

## 3.4 Algorithm

As this project focuses unsupervised, we mainly use two clustering algorithms Partitioning around medoids(PAM) and K prototype clustering algorithms for grouping the cars.

### i) Partitioning Around Medoids(PAM)

Step1: Randomly selects $k$ user models as $k$ medoids, each medoid represents a cluster

Step 2: Assigning user models into clusters according to similarity measure

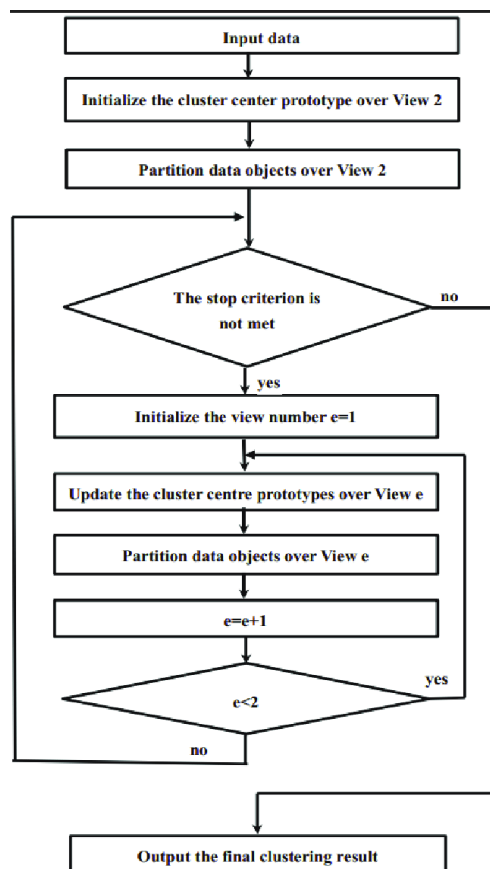Step3: Re-calculating k medoids according to similarity measure

No

Terminated criterion is satisfied?

Yes

Terminated

Partitioning Around Medoids (PAM) is a clustering algorithm that seeks to minimise the sum of dissimilarities between data points and their assigned medoids (representative points within each cluster). The steps of algorithm is listed below:-

1. Define the number of clusters (K) you want to form and randomly select K data points to be the initial medoids.

2. Assign each data point to the nearest medoid based on some distance metric (e.g. Euclidean distance, Manhattan distance, etc.).

3. For each cluster, calculate the total dissimilarity between each data point and its medoid.

4. For each cluster, select a non-medoid data point that would minimize the total dissimilarity of the cluster if it were to replace the current medoid.

5. Calculate the total dissimilarity of the new set of medoids.

6. If the total dissimilarity is lower than the previous set of medoids, update the medoids and go back to step 2. If the total dissimilarity is not lower, stop the algorithm and return the final set of medoids and their assigned data points.

7. Repeat steps 2-6 until the algorithm converges (i.e. the medoids no longer change).

ii) **K Prototype Clustering**

K-Prototype Clustering is a variant of K-means algorithm that can handle mixed data types by using a combination of categorical and numerical distance metrics. Here's a step-by-step algorithm for K-Prototype Clustering:

1. Define the number of clusters (K) you want to form and randomly select K data points to be the initial centroids.

2. Assign each data point to the nearest centroid based on the combination of categorical and numerical distance metrics. The categorical distance metric is usually a dissimilarity measure such as Hamming distance, while the numerical distance metric is typically the Euclidean distance.

3. For each cluster, calculate the average dissimilarity between each data point and its centroid. The average dissimilarity is the sum of the categorical dissimilarity and numerical distance divided by the number of data points in the cluster.

4. For each centroid, select a data point that would minimize the average dissimilarity of the cluster if it were to replace the current centroid. The selection process is based on a combination of the categorical and numerical distances.

5. Calculate the total average dissimilarity of the new set of centroids.

6. If the total average dissimilarity is lower than the previous set of centroids, update the centroids and go back to step 2. If the total average dissimilarity is not lower, stop the algorithm and return the final centroids and their assigned data points.

7. Repeat steps 2-6 until the algorithm converges (i.e. the centroids no longer change).
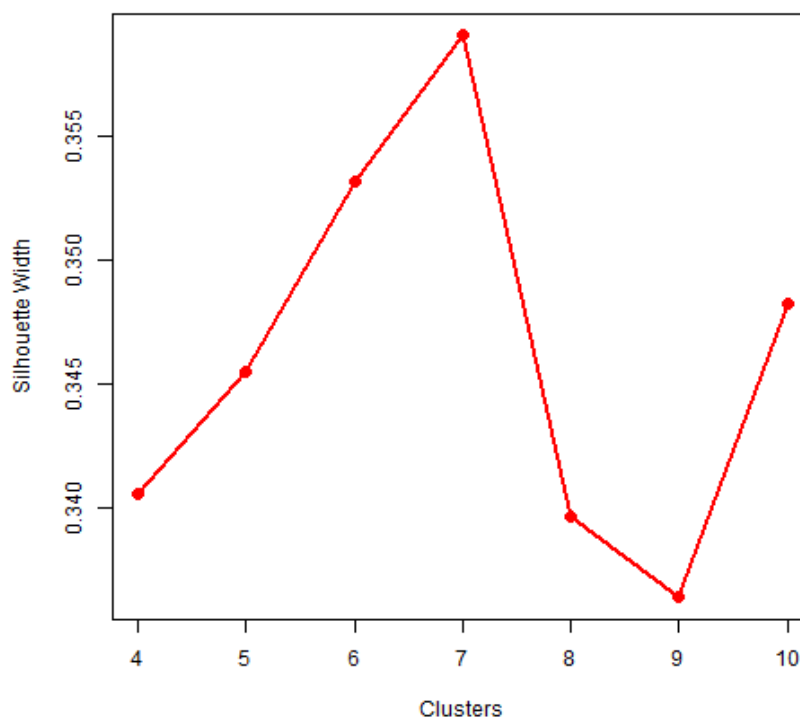
# Chapter 4 - Results And Discussions

## 4.1 Result & Output

### 4.1.1 Partitioning Around Medoids

The dataset to be clustered had both numerical and categorical data. Hence Partitioning Around Medoids, a.k.a the K-Medoids algorithm was used with Gower's Distance as the distance metric to cluster the car models. The results of the K-Medoids clustering was a model-wise segmentation of car models from multiple brands. By this, a car manufacturer would get to know a competitive car model of another brand for a particular car model.

I. **Silhouette Graph -**
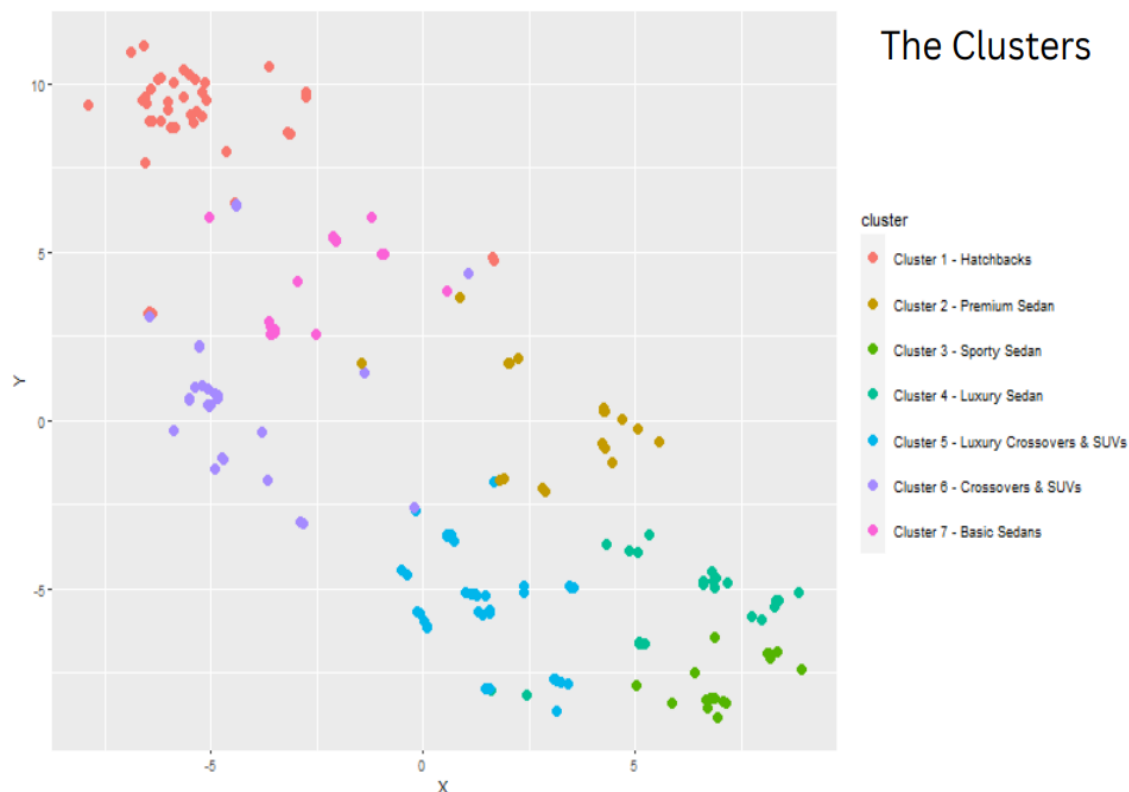
A silhouette graph was plotted with the number of clusters on the x-axis and the silhouette width on the y-axis to find out the optimal number of clusters. It is a measure of the quality of clustering for various number of clusters. The silhouette width ranges between -1 and 1 and higher the value, better the clustering. A peak was observed at K=7 and hence having 7 clusters would be optimum.

## II.    Cluster Visualisation -

The t-Distributed Neighbor Embedding (t-SNE) a statistical method that is mainly used to visualize high-dimensional data was used to plot the clusters.

### The Clusters

cluster

- Cluster 1 - Hatchbacks
- Cluster 2 - Premium Sedan
- Cluster 3 - Sporty Sedan
- Cluster 4 - Luxury Sedan
- Cluster 5 - Luxury Crossovers & SUVs
- Cluster 6 - Crossovers & SUVs
- Cluster 7 - Basic Sedans

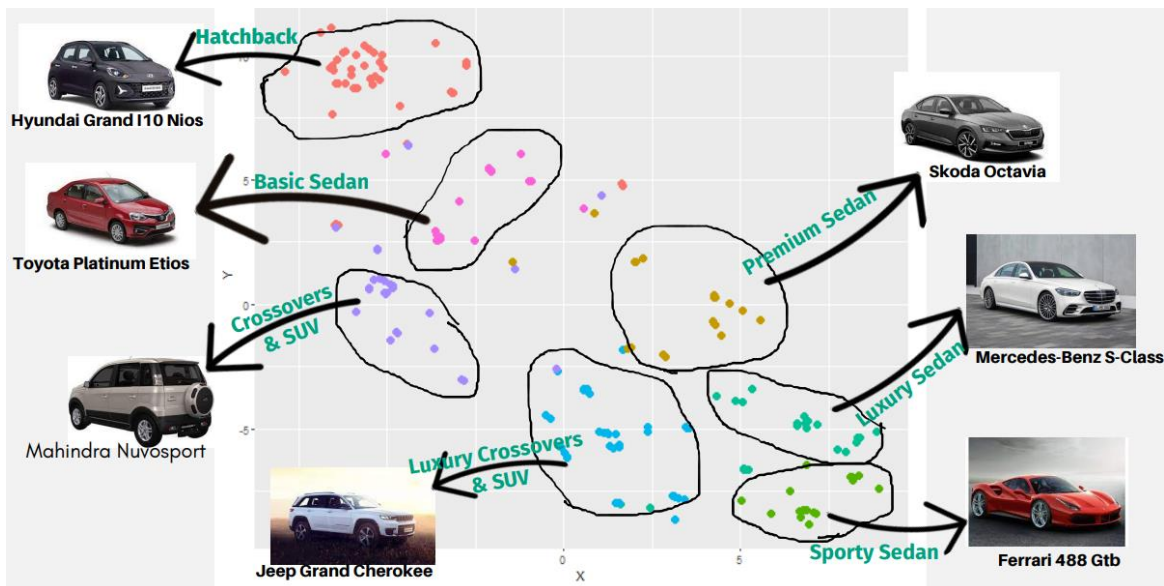## III.    Cluster Representatives

## Medoids

The Seven Clusters were named based on the car's price range, luxury level and body type. The cluster names are -

1. Cluster 1 - Hatchback
2. Cluster 2- Premium Sedan
3. Cluster 3 - Sporty Sedan
4. Cluster 4 - Luxury Sedan
5. Cluster 5 - Luxury Crossovers & SUV
6. Cluster 6 - Crossovers & SUV
7. Cluster 7 - Basic Sedan

The cluster representatives - the medoid car model of each cluster are tabulated below.

Description: df [7 × 3]

| car | Clusters | labels |
|-----|----------|--------|
| Hyundai Grand I10 Nios | 1 | Hatchbacks |
| Skoda Octavia | 2 | Premium Sedan |
| Ferrari 488 Gtb | 3 | Sporty Sedan |
| Mercedes-Benz S-Class | 4 | Luxury Sedan |
| Jeep Grand Cherokee | 5 | Luxury Crossovers & SUVs |
| Mahindra Nuvosport | 6 | Crossovers & SUVs |
| Toyota Platinum Etios | 7 | Basic Sedans |



The samples from each cluster are tabulated below -

## Cluster -1 (Hatchbacks)

Cluster -1 consists of Hatchback car models. A total of 45 car models belong to this cluster.

A tibble: 45 × 31

| car | Ex.Showroom_Price | Displacement | ARAI_Certified_Mileage |
|-----|-------------------|--------------|------------------------|
| Go+ | 571355.0 | 1198 | 26.07595 |
| Datsun Go | 512183.2 | 1198 | 25.44304 |
| Datsun Redi-Go | 358994.3 | 799 | 30.83966 |
| Fiat Punto Evo | 651820.8 | 1172 | 23.59816 |
| Fiat Punto Evo Pure | 529138.5 | 1172 | 21.51290 |
| Ford Ecosport | 985725.0 | 1497 | 25.61536 |
| Ford Figo | 645933.3 | 1498 | 28.23591 |
| Ford Freestyle | 725025.0 | 1498 | 26.68878 |
| Honda Jazz | 861187.5 | 1199 | 26.95659 |
| Hyundai Elite I20 | 779543.6 | 1197 | 24.97477 |

## Cluster - 2 (Premium Sedans)

Cluster -2 consists of Sedans from a premium range. Their prices are higher than the Basic Sedans (Cluster - 7). A total of 21 car models belong to this cluster.

| car | Ex.Showroom_Price | Displacement | ARAI_Certified_Mileage |
|---|---|---|---|
| Mercedes-Benz A-Class | 2836080 | 1595 | 19.12503 |
| Mercedes-Benz B-Class | 3145027 | 1595 | 18.38784 |
| Mercedes-Benz Cla-Class | 3490000 | 2143 | 20.68894 |
| Audi A3 | 3070975 | 1968 | 24.39949 |
| Audi A3 Cabriolet | 5038000 | 1395 | 24.30380 |
| Audi A6 | 5692200 | 1984 | 17.86076 |
| Honda Accord Hybrid | 4321237 | 1993 | 23.10000 |
| Honda Civic | 2030700 | 1799 | 25.49221 |
| Maruti Suzuki Baleno Rs | 788913 | 998 | 27.08861 |
| Mini Convertible | 3830000 | 1998 | 18.46835 |

## Cluster - 3 (Sporty Sedans)

Cluster - 3 consists of Sporty Sedans with higher price than the Basic Sedans (Cluster - 7). Sporty Sedans are characterised by a lower height, higher number of cylinders, usually 8 or 12 and upto 8 gears. A total of 19 car models belong to this cluster.

| car | Ex.Showroom_Price | Displacement | ARAI_Certified_Mileage |
|---|---|---|---|
| Mercedes-Benz Amg-Gt | 22527075 | 3982 | 9.873418 |
| Rolls-Royce Phantom Coupe | 77312661 | 6749 | 28.860759 |
| Rolls-Royce Wraith | 50025840 | 6592 | 12.911392 |
| Aston Martin Vantage | 29500000 | 3982 | 10.126582 |
| Audi R8 | 27245000 | 5204 | 8.493671 |
| Audi Rs7 | 17072000 | 3993 | 12.405063 |
| Bentley Continental Gt | 33791059 | 5998 | 10.886076 |
| Bmw M4 | 13590000 | 2979 | 13.607595 |
| Dc Avanti | 3407407 | 2000 | 12.658228 |
| Ferrari 458 Speciale | 42531500 | 4497 | 11.392405 |

## Cluster - 4 (Luxury Sedans)

Cluster -4 consists of high end luxury sedan models from top brands that are known for their luxury and lavishness. Their prices are much higher than the Basic Sedans (Cluster - 7). A total of 25 car models belong to this cluster.

| car | Ex.Showroom_Price | Displacement | ARAI_Certified_Mileage |
|---|---|---|---|
| Mercedes-Benz C-Class | 4984792 | 1950 | 318.831685 |
| Mercedes-Benz E-Class | 7263839 | 2987 | 92.377678 |
| Mercedes-Benz Glc | 5525000 | 1991 | 19.476114 |
| Mercedes-Benz Maybach | 23599156 | 5980 | 9.873418 |
| Mercedes-Benz S-Class | 17678364 | 2987 | 14.606966 |
| Rolls-Royce Dawn | 59216193 | 6598 | 10.000000 |
| Rolls-Royce Drophead Coupe | 83755383 | 6749 | 12.405063 |
| Rolls-Royce Ghost Series Ii | 49978467 | 6593 | 15.563291 |
| Rolls-Royce Phantom | 85200000 | 6749 | 29.113924 |
| Aston Martin Rapide | 38845823 | 5935 | 16.329114 |

## Cluster - 5 (Luxury Crossovers & SUVs)

Cluster -5 consists of crossovers and SUVs from top brands that are known for luxury. Their prices are higher than the Basic Crossovers & SUVs (Cluster - 6). Crossovers and SUVs are distinguished from other car types by a higher height, displacement and drive-train (4 wheel drive). A total of 38 car models belong to this cluster.

| car | Ex.Showroom_Price | Displacement | ARAI_Certified_Mileage › |
|---|---|---|---|
| Mercedes-Benz Gla-Class | 4991824 | 1991 | 19.491389 |
| Mercedes-Benz Gls | 10835216 | 2987 | 14.690847 |
| Audi Q3 | 3881788 | 1968 | 19.527789 |
| Audi Q5 | 5321200 | 1968 | 17.678105 |
| Audi Q7 | 7601500 | 2967 | 16.992477 |
| Bentley Bentayga | 41741813 | 5950 | 17.575949 |
| Bmw 6-Series | 6926667 | 2993 | 21.907051 |
| Bmw 7-Series | 15543333 | 2993 | 17.753298 |
| Bmw X1 | 4018000 | 1995 | 23.138602 |
| Bmw X3 | 5803333 | 1995 | 22.703019 |

## Cluster - 6 (Crossovers & SUVs)

Cluster - 6 consists of Crossovers & SUVs that are cheaper in price. They are characterised by a larger displacement, height and weight than the hatchbacks.  A total of 28 car models belong to this cluster.

| car | Ex.Showroom_Price | Displacement | ARAI_Certified_Mileage › |
|---|---|---|---|
| Fiat Avventura | 794758.0 | 1248 | 24.63942 |
| Force Gurkha | 1150500.0 | 2596 | 20.43269 |
| Honda Brv | 1197025.0 | 1497 | 22.14927 |
| Honda Wr-V | 933577.8 | 1498 | 26.87253 |
| Hyundai Creta | 1309603.4 | 1591 | 21.72363 |
| Hyundai Tucson | 2281189.9 | 1995 | 18.35910 |
| Hyundai Venue | 919961.5 | 998 | 25.01831 |
| Isuzu Mu-X | 2833184.0 | 2999 | 16.58654 |
| Mahindra Bolero | 861919.8 | 2523 | 19.18269 |
| Mahindra Bolero Power Plus | 818333.0 | 1493 | 19.70192 |

## Cluster - 7 (Basic Sedans)

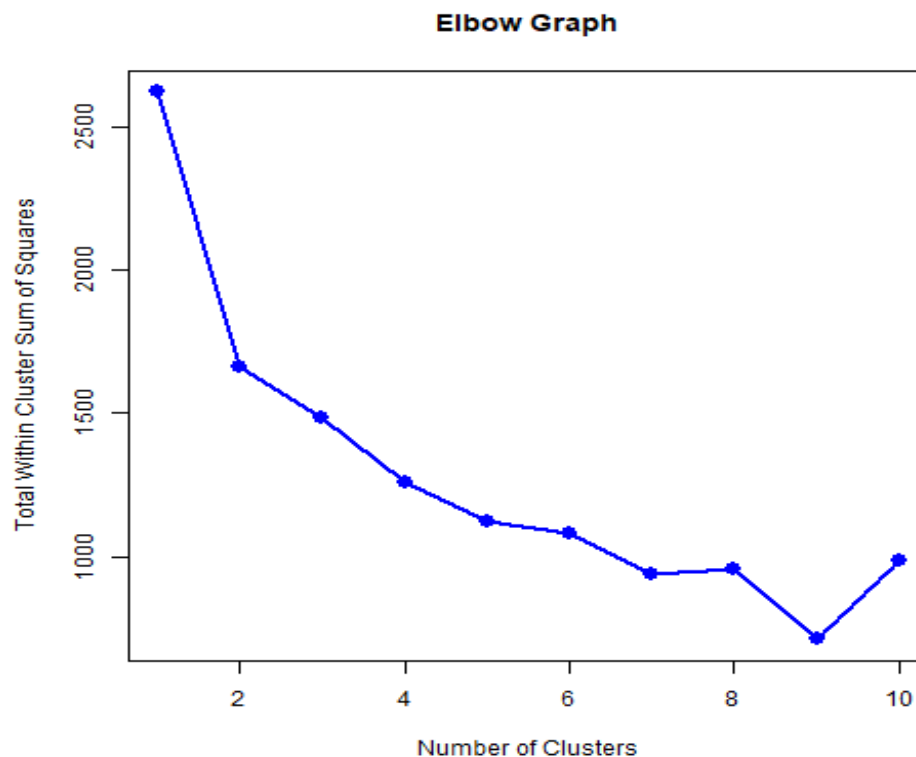Cluster - 7 consists of Basic Sedans. A total of 23 car models belong to this cluster.

| car | Ex.Showroom_Price | Displacement | ARAI_Certified_Mileage |
|-----|-------------------|--------------|------------------------|
| Fiat Linea | 899851.4 | 1248 | 22.66091 |
| Fiat Linea Classic | 735365.7 | 1248 | 21.91192 |
| Ford Aspire | 748939.5 | 1194 | 27.67388 |
| Honda Amaze | 792400.0 | 1199 | 27.07476 |
| Honda City | 1220790.0 | 1497 | 26.23813 |
| Hyundai Elantra | 1856500.0 | 1999 | 18.48101 |
| Hyundai Verna | 1131332.6 | 1591 | 23.49606 |
| Hyundai Xcent | 730617.8 | 1120 | 24.32068 |
| Hyundai Xcent Prime | 574315.2 | 1197 | 23.62372 |
| Mahindra Verito | 788713.7 | 1461 | 25.24038 |

## 4.1.2 K-Protoype

The K-Prototype algorithm provides segmentation based on brands. All the luxurious and extravagant brands are grouped together while the basic and economic brands fall into a group. The automobile market segmentation is at the brand level. A manufacturer can get to know their competitors in the Indian Market.
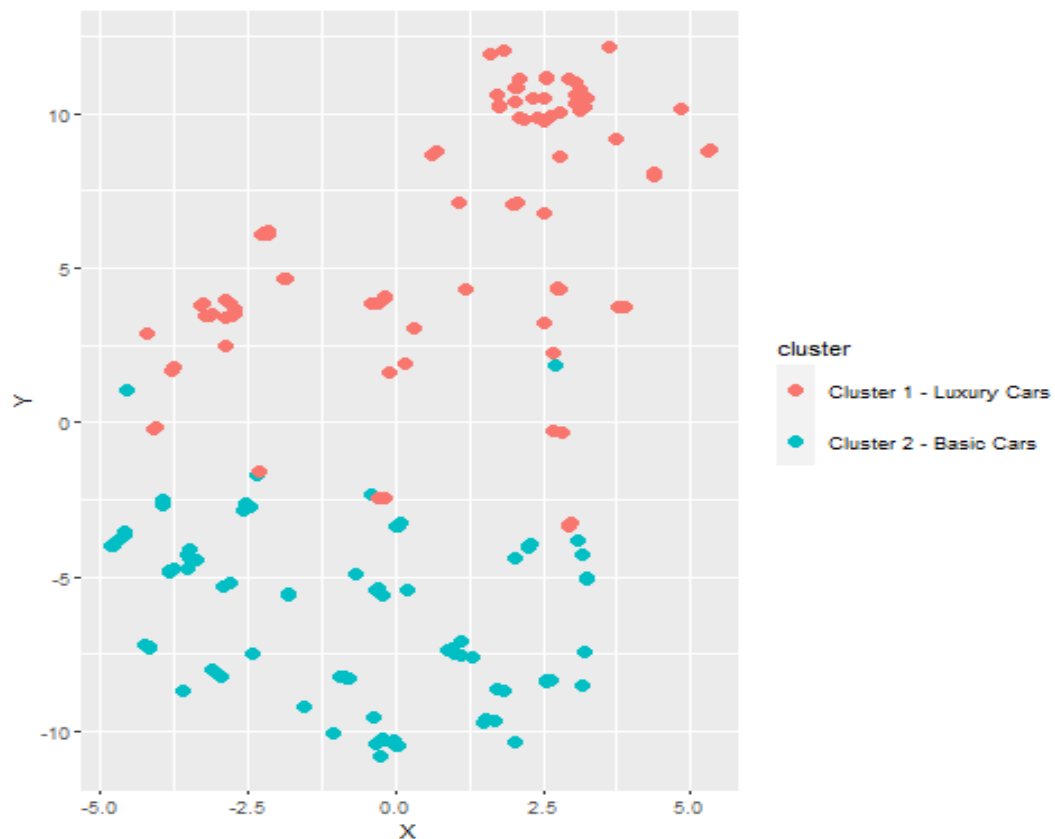
I. **Elbow Graph -**

An elbow graph with Number of Clusters on the x-axis and the Total Within-Cluster Sum of Squares on the y-axis was plotted to find the optimal number of clusters. The elbow bend occurs for Number of Clusters $= 2$ and hence the optimal number of clusters for the K-prototype model is fixed to be 2.



Elbow Graph

## II.  Cluster Visualisation -

The t-Distributed Neighbor Embedding (t-SNE) a statistical method that is mainly used to visualize high-dimensional data was used to plot the clusters.  The Cluster-1 represents Luxury Cars while the Cluster-2 represents economic cars.



## III.  Cluster Representatives
## Centroids

The k-Prototype algorithm computes cluster centroids rather than assigning a data instance as medoid. The cluster centroids are -

| Drivetrain | Fuel_Type | Body_Type | Rear_Brakes | Seats_Material |
|---|---|---|---|---|
| AWD (All Wheel Drive) | Petrol | SUV | Ventilated Disc | Leather |
| FWD (Front Wheel Drive) | Petrol | Hatchback | Drum | Fabric |

## Cluster - 1 (Basic Brands)

Cluster - 1 consists of car models from Economic Brands. This group consists of multiple body-types, and are characterized by their cheaper price. A total of 102 car models belong to this cluster. A few of the car models belonging to this cluster are-

| car | Ex.Showroom_Price | Displacement | ARAI_Certified_Mileage › |
|-----|-------------------|--------------|--------------------------|
| Ford Aspire | 748939.5 | 1194 | 27.67388 |
| Ford Ecosport | 985725.0 | 1497 | 25.61536 |
| Ford Figo | 645933.3 | 1498 | 28.23591 |
| Ford Freestyle | 725025.0 | 1498 | 26.68878 |
| Honda Accord Hybrid | 4321237.0 | 1993 | 23.10000 |
| Honda Amaze | 792400.0 | 1199 | 27.07476 |
| Honda Brv | 1197025.0 | 1497 | 22.14927 |
| Honda City | 1220790.0 | 1497 | 26.23813 |
| Honda Civic | 2030700.0 | 1799 | 25.49221 |
| Honda Jazz | 861187.5 | 1199 | 26.95659 |
| Hyundai Venue | 919961.5 | 998 | 25.01831 |
| Hyundai Verna | 1131332.6 | 1591 | 23.49606 |
| Hyundai Xcent | 730617.8 | 1120 | 24.32068 |
| Hyundai Xcent Prime | 574315.2 | 1197 | 23.62372 |
| Mahindra Bolero | 861919.8 | 2523 | 19.18269 |
| Mahindra Bolero Power Plus | 818333.0 | 1493 | 19.70192 |
| Mahindra Kuv100 Nxt | 660978.4 | 1198 | 25.00869 |
| Mahindra Nuvosport | 923606.5 | 1493 | 20.28446 |
| Mahindra Scorpio | 1321884.3 | 2179 | 17.56410 |
| Mahindra Thar | 977956.0 | 2498 | 15.62500 |

## Cluster - 2 (Luxury Brands)

Cluster - 2 consists of car models from Luxury Brands. This group consists of multiple body-types, and are characterized by their high price. A total of 97 car models belong to this cluster. A few of the car models belonging to this cluster are-

| car | Ex.Showroom_Price | Displacement | ARAI_Certified_Mileage › |
|-----|-------------------|--------------|--------------------------|
| Mercedes-Benz Amg-Gt | 22527075 | 3982 | 9.873418 |
| Mercedes-Benz B-Class | 3145027 | 1595 | 18.387841 |
| Mercedes-Benz C-Class | 4984792 | 1950 | 318.831685 |
| Mercedes-Benz Cla-Class | 3490000 | 2143 | 20.688940 |
| Mercedes-Benz E-Class | 7263839 | 2987 | 92.377678 |
| Mercedes-Benz Gla-Class | 4991824 | 1991 | 19.491389 |
| Mercedes-Benz Glc | 5525000 | 1991 | 19.476114 |
| Mercedes-Benz Gls | 10835216 | 2987 | 14.690847 |
| Mercedes-Benz Maybach | 23599156 | 5980 | 9.873418 |
| Mercedes-Benz S-Class | 17678364 | 2987 | 14.606966 |
| Audi A6 | 5692200 | 1984 | 17.860759 |
| Audi Q3 | 3881788 | 1968 | 19.527789 |
| Audi Q5 | 5321200 | 1968 | 17.678105 |
| Audi Q7 | 7601500 | 2967 | 16.992477 |
| Audi R8 | 27245000 | 5204 | 8.493671 |
| Audi Rs7 | 17072000 | 3993 | 12.405063 |
| Audi S5 | 7243000 | 2995 | 17.177215 |
| Bentley Bentayga | 41741813 | 5950 | 17.575949 |
| Bentley Continental Gt | 33791059 | 5998 | 10.886076 |
| Bentley Flying Spur | 36774574 | 3993 | 5.063291 |

### 4.1.3 Performance Evaluation

The cluster statistics for both the K-Medoids and K-Prototype algorithm are summarised below -

| Metric | K - Medoid | K - Prototype |
|---|---|---|
| Number of Clusters | 7 | 2 |
| Average between cluster distance | 0.4118 | 0.4875 |
| Average within cluster distance | 0.1596 | 0.2566 |
| Maximum cluster diameter | 0.5500 | 0.6523 |
| Minimum Separation | 0.0066 | 0.0124 |
| Average Silhouette width | 0.3590 | 0.4557 |
| Dunn Index | 0.0121 | 0.0190 |
| Calinski Harabasz Index | 120.184 | 208.517 |

Based on the Calinski Harabasz Index and the Dunn Index, K-Prototype algorithm seems to perform well than the K-Medoids algorithm with minimum between cluster similarities. Higher values of Calinski Harabasz Index and the Dunn Index support this.

## 4.2 Inference

### 4.2.1 Inferences from K-Medoids Clustering

#### 1) Brands Positions

Maruti Suzuki, Tata, Mahindra, Hyundai have clearly positioned themselves only in the economic and basic segment, largely targeting middle class consumers. It can be clearly seen from the clustering results that the above Makers occur only in the Basic Sedans, Hatchbacks and thus they are strong competitors of each other when it comes to the economic market. Likewise, Mercedes, BMW, Audi, Porsche and Jaguar compete with each other in the Luxury segment. They are prime manufacturers of sporty and luxury sedans.

#### 2) Honda City vs Honda Amaze

Honda City is a higher priced and premium sedan than Honda Amaze, but that does not seem true from the clustering results. Both car models are grouped into the same segment, which may not be the desired target market perceived by Honda. This suggests that Honda City being a premium car model than a more basic Honda Amaze, may often be perceived alike by most consumers.

#### 3) Unpopularity of Isuzu Mu-X

Isuzu, a Japanese car manufacturer, did not make much in the Indian automobile market. From 2013 until March 2022, only 2598 units of the car were sold. This could possibly be because of tough competitors in its cluster like Toyota Innova Crysta and Hyundai Creta which are highly sought after by buyers. Isuzu ' s failure may be because of their choice of target market - the SUV segment ruled by well established brands like Toyota and Hyundai.

### 4) Volkswagen Ameo is a perfect successor to Volkswagen Vento

Volkswagen wanted to introduce Ameo as a replacement to Volkswagen Vento, which seems successful from the clustering results. Both the car models occur in the low end sedan segment and hence indicates how Volkswagen continues its hold in the low end sedan segment.

### 5) Toyota Fortuner vs Ford Endeavour

The clustering result suggests a tough competition between Toyota Fortuner and Ford Endeavour. Both occurring in the Luxury Crossovers & SUVs segment, Ford's Endeavour competes with Toyota ' s Fortuner in every aspect including price, displacement and mileage.

### 6) Skoda Superb's Competitors

Skoda Superb competes with Volkswagen Passat and even the basic models from Mercedes and Audi in the luxury sedan segment. Because of an equal competition in the market between the cars in this segment, none outstands the other. This segment is equally held by makers like Mercedes, Volkswagen, Skoda and Audi.

### 7) Hyundai Creta vs Hyundai Tucson

Hyundai is a competitor to itself when it comes to Creta vs Tucson. The newly launched Tucson is not as sought after as Creta. Annual sales of Creta was 1,40,895 in 2022 while Tuscon was only a few thousands. Both occur in the Crossover & SUV segment and hence the same target market. This means that trying to increase the sales of one model will bring down the sales of the other.

### 8) Maruti Suzuki's Monopoly

Clearly seen from the Clustering results, Maruti Suzuki dominates the Hatchback segment with a total of 15 car models and contributes a 33% market share in this strategic group. Maruti Suzuki thus holds a major share in the Indian automobile market, primarily dominating the hatchback segment.

### 4.2.2 Inferences from K-Prototype Clustering

### 1) Brand Segmentation

The K-Prototype Clustering provides a segmentation of brands. It results in grouping the Indian Automobile Brands into two strategic groups -

1. Luxury Brands - The luxury brands consist of Mercedes, Audi, Porsche, Rolls Royes to name a few. They are leading manufacturers of luxury cars across the globe. These brands are sought after by the well off consumers.

2. Economic Brands - Brands like Maruti, Mahindra, Hyundai, Honda, Fiat, Ford, etc. fall under this segment. Their cars are economical and are widely sought after by the Indian automobile consumer market. They dominate the automobile market in terms of their sales. The consumers of these brands belong to middle and upper middle class.

## 4.3 Performance Analysis

K-medoids tends to have higher computational complexity compared to k-prototype because it involves finding the best medoid from a set of data points in each iteration, which can be computationally expensive, especially for large datasets. K-prototype, on the other hand, involves updating centroids based on numerical means and categorical modes, which can be computationally more efficient.

- The time complexity of the k-prototypes algorithm :  **$O(k*p*n)$**
- The time complexity of K-medoids algorithm is : **$O(k * (n-k)^2)$**

where 'n' is the number of data points, 'k' is the number of clusters and 'p' is the number of features

# Chapter 5: Conclusion

The "Indian Automotive Strategic Grouping" study represents a valuable contribution to the field of Indian automotive market analysis. With its utilisation of a data-driven approach and a clustering method, the study provides a comprehensive and nuanced understanding of consumer preferences and buying behaviour. By grouping automobile models into distinct strategic groups based on various attributes, including price, brand, segment, and features, the study offers valuable insights for automobile manufacturers, dealers, and marketers seeking to develop targeted marketing strategies and improve product offerings.Customers can use this information to make informed purchasing decisions, identifying the car models that best match their preferred attributes and provide the best value for money.Additionally, the study identifies the best car model in each cluster based on value for money and features, providing customers with a clear idea of which models offer the most desirable features for the best price. Overall, this study can help customers make informed purchasing decisions, by providing valuable insights into the Indian automotive market, and identifying the car models that best match their needs and preferences. As such, the study is highly recommended for anyone with an interest in the Indian automotive market or automotive market analysis more broadly.

# Chapter 6: Future Scope

The "Indian Automotive Strategic Grouping" study offers valuable insights into the Indian automotive market, providing a comprehensive and nuanced understanding of consumer preferences and buying behaviour. As such, it has several future scopes that can enhance its value and contribute to the field of Indian automotive market analysis. Some of the feature scope has been discussed below:-

1. The study could be updated regularly to reflect changes in the Indian automotive market, such as the introduction of new car models or changes in consumer preferences. This would ensure that the study remains relevant and up-to-date, providing valuable insights for industry professionals and customers alike.
2. This could also be expanded to include a more comprehensive analysis of the Indian automotive market, covering aspects such as sales trends, market share, and competitive landscape. This would provide a more complete understanding of the Indian automotive market, helping industry professionals make informed decisions and enhancing the study's overall value.
3. Finally, the study's insights could be applied to develop new products and services that cater to the needs and preferences of Indian consumers. For instance, manufacturers could use the study's clustering method to develop targeted marketing strategies or new car models that offer features and attributes that are most valued by Indian consumers. This could help manufacturers improve their market position and competitiveness in the Indian automotive market.

Overall, the "Indian Automotive Strategic Grouping" study has several future scopes that can enhance its value and contribute to the field of Indian automotive market analysis. Its insights can help industry professionals and customers make informed decisions, and its application can lead to the development of new products and services that better cater to Indian consumers' needs and preferences.

# Chapter 7: References

1. K. Noor and S. Jan, "Vehicle price prediction system using machine learning techniques, " International Journal of Computer Applications, vol. 167, no. 9, pp. 27–31, 2017.
2. Mou, A. D., Saha, P. K., Nisher, S. A., & Saha, A. (2021, February). A comprehensive study of machine learning algorithms for predicting car purchase based on customers demands. In 2021 International Conference on Information and Communication Technology for Sustainable Development (ICICT4SD) (pp. 180-184). IEEE.
3. Dave, M., & Gianey, H. (2016, November). Different clustering algorithms for Big Data analytics: A review. In 2016 International Conference System Modeling & Advancement in Research Trends (SMART) (pp. 328-333). IEEE
4. Porter, M. E. (1980). The strategic grouping of automobile firms: A review of literature. Academy of Management Review, 5(4), 519-527.
5. Chen, Y. S., & Chen, Y. J. (2013). Strategic grouping and competitive dynamics in the global automobile industry. Journal of Business Research, 66(4), 493-500.
6. Cutler, D. M., & Poterba, J. M. (1989). Strategic grouping in the US automobile industry. Journal of Economics & Management Strategy, 7(2), 143-166.
7. Aoki, M., & Hayami, Y. (1986). The effect of strategic grouping on firm performance: Evidence from the Japanese automobile industry. Journal of Industrial Economics, 35(3), 245-256.
8. Qin, X., & Liu, J. (2018). Strategic grouping in the Chinese automobile industry. Journal of Business Research, 87, 99-108.
9. Fri, W., Pehrsson, T., & Søilen, K. (2013). How phases of cluster development are associated with innovation-the case of China. International journal of innovation science, 5(1), 31-44.
10. Okada, A., & Siddharthan, N. S. (2007). Industrial clusters in India: Evidence from automobile clusters in Chennai and the national capital region.

11. Moustaki, I., & Knott, M. (2000). A latent variable model for mixed mode data: The Mixture-of-Distribution approach. Journal of Computational and Graphical Statistics, 9(2), 328-347.

12. Bhattacharyya, D., & Mukhopadhyay, A. (2015). Cluster validation using Calinski–Barabasz index for gene expression data. Journal of Statistical Planning and Inference, 157, 45-55.

13. Huang, Z. (1998). Extensions to the k-means algorithm for clustering large data sets with categorical values. Data Mining and Knowledge Discovery, 2(3), 283-304.

14. Kaufman, L., & Rousseeuw, P. J. (1990). Finding groups in data: An introduction to cluster analysis (Vol. 344). John Wiley & Sons