# ELECTRIC VEHICLE MARKET SEGMENT ANALYSIS



**"Electric cars are the future." – Albert II, Prince of Monaco.**

SUBMITTED BY SAUMYA P JOHN

22nd August 2024

**GitHub Link:**

https://github.com/Saumyapj/FeynnLabsInternship/blob/main/Saumya_John_EVsegmentation_Analysis.ipynb

**Abstract:**

The electric vehicle (EV) market in India is experiencing rapid growth, driven by increasing environmental awareness, government incentives, and advancements in technology. This study aims to analyze the EV market in India through segmentation analysis, focusing on various factors such as vehicle type, geographic distribution, and consumer demographics. Emphasizing interpretability, the approach provides actionable insights for tailored marketing campaigns and product strategies. This enables the startup to effectively target and retain customers in the competitive EV market.

**Methodology**: The analysis employs a combination of clustering techniques and principal component analysis (PCA) to identify distinct market segments. Data is collected from various sources, including government reports, industry publications, and consumer surveys. Key variables include vehicle types (two-wheelers, three-wheelers, four-wheelers, goods vehicles, public service vehicles, special category vehicles, ambulances/hearses, construction equipment vehicles, and others), geographic regions (urban vs. rural, state-wise distribution), and consumer demographics (age, income levels, education levels).

**Findings**: The segmentation analysis reveals significant variations in EV adoption across different segments. Two-wheelers dominate the market, particularly in urban areas, due to their affordability and convenience. Four-wheelers are gaining traction among higher-income groups and in metropolitan regions. The availability of charging infrastructure is a critical factor influencing EV adoption, with states having more charging stations showing higher EV penetration rates.

**Implications**: The insights from this segmentation analysis can help policymakers, manufacturers, and investors make informed decisions. Policymakers can tailor incentives and infrastructure development to target specific segments. Manufacturers can develop marketing strategies and product offerings that cater to the preferences of different consumer groups. Investors can identify high-potential segments and regions for investment.

**Conclusion**: The EV market in India is diverse and dynamic, with distinct segments showing varying levels of adoption and preferences. Understanding these segments is crucial for driving the growth of the EV market and achieving sustainable transportation goals.

# 1. Introduction:

**Problem Statement:**

- The electric vehicle (EV) market is booming due to rising environmental concerns, government incentives, and technological advancements. However, increased competition necessitates differentiation for EV startups. Traditional demographic-based segmentation methods are inadequate for this diverse market. Hence, there's a rising demand for sophisticated segmentation strategies like machine learning.

- Machine learning analyzes vast datasets to reveal subtle patterns and relationships, allowing startups to identify distinct market segments based on demographics, psychographics, and behaviour. This enables tailored marketing campaigns and product offerings. Moreover, the interpretability of machine learning models provides actionable insights, aiding in strategy development.

- Overall, applying machine learning for EV startup market segmentation provides a competitive edge by accurately identifying customer segments, enabling strategic positioning, effective engagement, and sustainable growth amidst fierce competition.

**Objectives of the project:**

- Identify Optimal EV Type: Determine the ideal location and electric vehicle (EV) model for launch by analyzing market trends, consumer preferences, and technological feasibility.

- Segment Customer Base: Utilize machine learning to identify distinct customer segments within the EV market based on demographics, behaviour, and psychographics.

# 2. Data Collection and Preprocessing:

**Data collection process & Description of the dataset used**

To kickstart our EV startup's market segmentation analysis for the upcoming launch in India, I began by focusing on data acquisition. This involved extensive research across multiple online sources to gather pertinent and suitable data for our project. This thorough data gathering process forms the foundation for the next pivotal phase: pinpointing the most lucrative segment to ensure a successful entry into India's dynamic and burgeoning EV market.

**Resources used for research:**

- https://www.kaggle.com/
- https://data.gov.in/

- https://cea.nic.in/electric-vehicle-charging-reports/?lang=en
- https://dataspace.mobi/dataset/electric-vehicle-charging-station-list
- https://www.statista.com/statistics/1264923/india-electric-passenger-vehicle-sales-by-manufacturers/
- https://vahan.parivahan.gov.in/vahan4dashboard/vahan/view/reportview.xhtml
- https://datasetsearch.research.google.com/

**EV_Customer_preferences.csv :** This dataset comprises information on customer demographics, preferences and perception. Dataset includes Age, City, Profession, Marital status, Income and their opinion and preferences.

**final_dataset.csv**: This dataset comprises information on state wise different vehicle types total number of vehicles and charging stations across India.

**EV_cars_India_2023.csv**: This dataset comprises information on car name and their speed, price, boot space, power and charging time

## 3. Data Pre-processing:(Steps and libraries used)

Importing Libraries: firstly, we will import the libraries for our model

```python
#importing the Libraries
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from scipy import stats
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn import preprocessing
from sklearn.cluster import KMeans
from collections import Counter
from itertools import product
from bioinfokit.visuz import cluster
from sklearn.cluster import KMeans
from yellowbrick.cluster import KElbowVisualizer
from statsmodels.graphics.mosaicplot import mosaic
import math
import warnings
warnings.filterwarnings('ignore')
```

- Numpy we have imported for the performing mathematics calculation.
- Matplotlib is for plotting the graph, and pandas are for managing the dataset.
- Seaborn is for data visualization library; it is based on matplotlib.
- Scikit-learn have sklearn.cluster.KMeans module to perform K-Means clustering. While
- computing cluster centers and value of inertia.

- The bioinfokit toolkit aimed to provide various easy-to-use functionalities to analyze,

  visualize, and interpret the biological data.

- The KElbowVisualizer implements the "elbow" method to help data scientists select the optimal number of clusters by fitting the model with a range of values for K.

- Mosaic plots help show relationships and give a visual way to compare groups.

**Fetching the dataset**

```python
# fetching EV Customer preference- dataset1
df = pd.read_csv('EV_Customer_preferences.csv')
```

**Checking the null values and fill the null values**

```python
print(df.isna().sum())
```

```python
df['Profession']=df['Profession'].fillna(df['Profession'].mode()[0])
```

**Checking for Duplicates**

```python
print(f"Duplicates: {df.duplicated().sum()}")
duplicate = df[df.duplicated()]
print("Duplicate Rows :")
# Print the resultant rows
duplicate
```

```
Duplicates: 0
Duplicate Rows :
```

**Checking for Outliers and handling them**

```python
#Boxplot for numerical column
for i in ['Age','Annual_Income','Preference for wheels in EV']:
    plt.figure()
    plt.boxplot(df[i])
    plt.title(i)
```

```
#Treating outliers
Q1 = df.Age.quantile(0.25)
Q3 = df.Age.quantile(0.75)
IQR = Q3 - Q1
min_limit = Q1 - (IQR * 1.5)
max_limit = Q3 + (IQR * 1.5)
print(min_limit, max_limit,IQR)
```

```
18.5 38.5 5.0
```

```
df.loc[df['Age']>max_limit,'Age']=np.mean(df.Age)
df.loc[df['Age']<min_limit,'Age']=np.mean(df.Age)
```

- **Encoding Categorical Variables**: LabelEncoder used to encode categorical variables.

```
# List of columns to encode
le=LabelEncoder()
columns_to_encode = [
    'City', 'Profession', 'Marital Status', 'Education',
    'Would you prefer replacing all your vehicles to Electronic vehicles?',
    'If Yes/Maybe what type of  EV would you prefer?',
    'Do you think Electronic Vehicles are economical?',
    'Which brand of vehicle do you currently own?',
    'How much money could you spend on an Electronic vehicle?',
    'Do you think Electronic vehicles will replace fuel cars in India?'
]
# Apply LabelEncoder to each column
for column in columns_to_encode:
    df[column] = le.fit_transform(df[column])

# Display the DataFrame
df
```

- **Feature Scaling**: StandardScaler() techniques applied to the selected features for modelling.

```
# feature scaling
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
```

## 4. Exploratory Data Analysis (EDA)

**Descriptive Statistics:**

```
df.describe()
```

|  | Age | No. of Family members | Annual_Income | Preference for wheels in EV |
|---|---|---|---|---|
| count | 1000.000000 | 1000.000000 | 1.000000e+03 | 1000.000000 |
| mean | 31.800000 | 4.118000 | 2.258342e+06 | 3.349000 |
| std | 11.294847 | 1.469774 | 9.993558e+05 | 0.887686 |
| min | 15.000000 | 0.000000 | -3.761509e+05 | 2.000000 |
| 25% | 26.000000 | 4.000000 | 1.782116e+06 | 2.000000 |
| 50% | 29.000000 | 4.000000 | 2.329246e+06 | 4.000000 |
| 75% | 31.000000 | 5.000000 | 2.753170e+06 | 4.000000 |
| max | 118.000000 | 8.000000 | 1.282128e+07 | 4.000000 |

**Visualization:**

**univariate analysis of categorical variables:**

.**Observations**

- Most of the professionals are working professional.
- Most of the peoples are Singles.
- Most of them are graduates.
- Most of them are prefer to replace their vehicles to EV and they like SUV.
- Most of the people prefer Tata brand.
- Most of the people prefer less than 15 lakhs vehicle.
- Most of the customers are from Pune city.



City preferences for EV

Most people are preferring 4 wheelers and it is more in Pune city.

**Bivariate analysis:**



Most of the middle-class income bracket people like to replace their vehicle to EV and ready to spend less than 15 lakhs for that.

## 5. Methodology and Machine Learning Modeling:

### Machine learning algorithms used:

K-means clustering was chosen due to its simplicity and effectiveness in segmenting data based on similarity. K-Means Clustering and Principal Component Analysis (PCA) can be used in tandem to analyze and extract insights from attributes.

PCA is a dimensionality reduction technique that can help reduce the complexity of geographic data while preserving essential information. It's particularly useful when dealing with high-dimensional data.

```python
# applying Principle Component Analysis (PCA)
pca = PCA(n_components=6)
X_pca = pca.fit_transform(X_scaled)
df_pca = pd.DataFrame(X_pca, columns=['PC1', 'PC2', 'PC3', 'PC4', 'PC5', 'PC6'])
df_pca.head()
```



**Number of clusters:** The optimal number of clusters was determined using the Elbow method.

```python
# plotting the results of Elbow

wcss = []

for i in range(1, 12):
    kmean = KMeans(n_clusters=i, init='k-means++', random_state=90)
    kmean.fit(X_pca)
    wcss.append(kmean.inertia_)

plt.figure(figsize=(8,6))
plt.title('Plot of the Elbow Method', size=15, family='serif')
plt.plot(range(1, 12), wcss)
plt.xticks(range(1, 12), family='serif')
plt.yticks(family='serif')
plt.xlabel('Number of Custers (K)', family='serif')
plt.ylabel('WCSS', family='serif')
plt.grid()
plt.tick_params(axis='both', direction='inout', length=6, color='purple', grid_color='lightgray', grid_linestyle='--')
plt.show()
```

Plot of the Elbow Method

Based on the elbow method, the optimal number of clusters for the given data is likely k=6. This means that the KMeans algorithm can effectively group the data points into six distinct clusters with minimal distortion.

```
# Applying K-means clustering
kmeans = KMeans(n_clusters=6,random_state=90)
df['Cluster'] = kmeans.fit_predict(X_scaled)

# analyze the clusters
print(df['Cluster'].value_counts())

Cluster
2    215
3    186
1    159
4    154
5    150
0    136
Name: count, dtype: int64
```

Interpreting Results:

```
features = X[['Age', 'City', 'Profession', 'Marital Status', 'Education', 'Annual_Income']]
kmeans = KMeans(n_clusters=6, random_state=90)
X['Segment'] = kmeans.fit_predict(features)
# Analyze the segments
print(X.groupby('Segment').mean())
```

|         | Age       | City      | Profession | Marital Status | Education |
|---------|-----------|-----------|------------|----------------|-----------|
| Segment |           |           |            |                |           |
| 0       | 27.992201 | 15.919414 | 1.560440   | 0.637363       | 0.274725  |
| 1       | 28.245600 | 15.569343 | 1.503650   | 0.627737       | 0.343066  |
| 2       | 28.480952 | 14.380952 | 1.619048   | 0.666667       | 0.333333  |
| 3       | 28.074781 | 15.570850 | 1.696356   | 0.668016       | 0.336032  |
| 4       | 28.690052 | 15.952880 | 1.727749   | 0.581152       | 0.272251  |
| 5       | 29.003636 | 14.936364 | 1.636364   | 0.672727       | 0.327273  |

|         | Annual_Income |
|---------|---------------|
| Segment |               |
| 0       | 2.521751e+06  |
| 1       | 1.705649e+06  |
| 2       | 7.025723e+05  |
| 3       | 2.972029e+06  |
| 4       | 2.108713e+06  |
| 5       | 1.215123e+06  |

```
segment = replace.merge(prefer, on='Cluster Number', how='left').merge(age, on='Cluster Number', how='left')
print(segment)
plt.figure(figsize = (12,4))
sns.scatterplot(x="Age", y="If Yes/Maybe what type of  EV would you prefer?", data=segment, hue='Cluster Number',s=400, palette='coolwarm')
plt.title("Simple segment evaluation plot EV", fontsize=15)
plt.xlabel("Age", fontsize=10)
plt.ylabel("What type of  EV would you prefer?", fontsize=10)
plt.legend(title='Segment')
plt.show()
```

```
   Cluster Number  \
0               0
1               1
2               2
3               3
4               4
5               5

   Would you prefer replacing all your vehicles to Electronic vehicles?  \
0                                           1.597059
1                                           1.450119
2                                           1.046512
3                                           1.510753
4                                           1.448052
5                                           1.500000

   If Yes/Maybe what type of  EV would you prefer?       Age
0                                       2.860294  24.970588
1                                       2.690113  28.573585
2                                       2.609302  28.198917
3                                       2.833333  30.401456
4                                       2.733766  28.569138
5                                       2.766667  28.385491
```



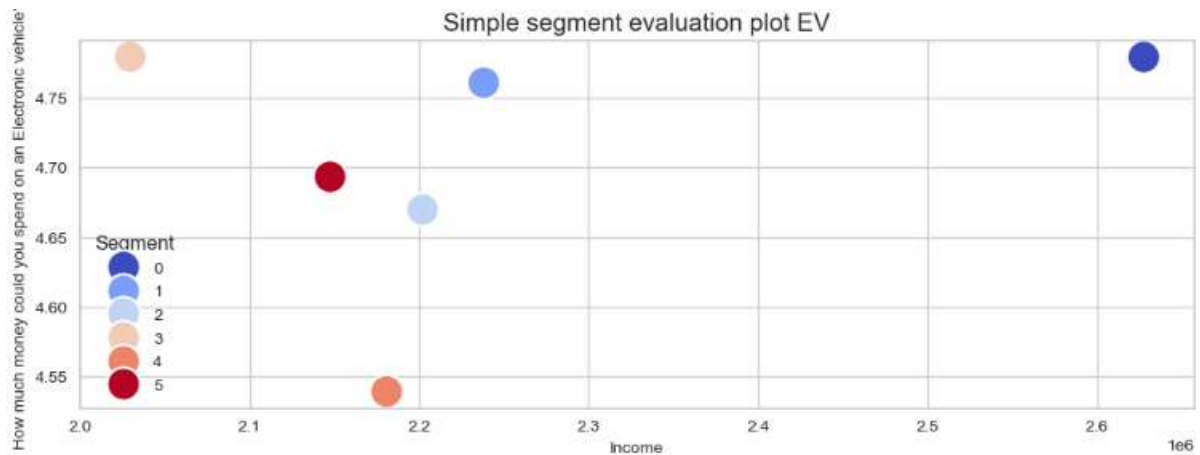Simple segment evaluation plot EV

### Observations:

The data points are concentrated in 4th part of the graph, suggesting that there is a unity in terms of both type of EVs and prefer to replacing the EVs across the identified segments. Most of the age groups prefer to buy SUVs compared to other models when they like to replace their vehicles to EV.

```
segment = replace.merge(income, on='Cluster Number', how='left').merge(spend, on='Cluster Number', how='left')
print(segment)
plt.figure(figsize = (12,4))
sns.scatterplot(x="Annual_Income", y="How much money could you spend on an Electronic vehicle?", data=segment, hue='Cluster Number',s=400, palette='cool
plt.title("Simple segment evaluation plot EV", fontsize=15)
plt.xlabel("Income", fontsize=10)
plt.ylabel("How much money could you spend on an Electronic vehicle?", fontsize=10)
plt.legend(title='Segment')
plt.show()
```

```
   Cluster Number  \
0               0
1               1
2               2
3               3
4               4
5               5

   Would you prefer replacing all your vehicles to Electronic vehicles?  \
0                                           1.597059
1                                           1.450119
2                                           1.046512
3                                           1.510753
4                                           1.448052
5                                           1.500000

   Annual_Income  How much money could you spend on an Electronic vehicle?
0   2.627455e+06                                          4.770412
1   2.230055e+06                                          4.761806
2   2.202125e+06                                          4.665767
3   2.020655e+06                                          4.773570
4   2.100055e+06                                          4.598061
5   2.167055e+06                                          4.803111
```
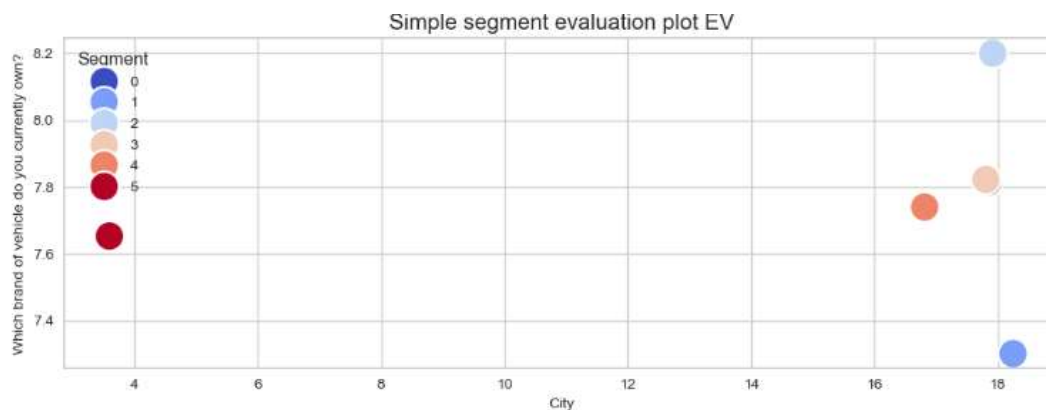
Simple segment evaluation plot EV

**Observations:**

The data points are concentrated in 2nd part of the graph, most middle-class income people are ready to spend >15 lakhs for their next EV.

```
segment = economical.merge(city, on='Cluster Number', how='left').merge(own, on='Cluster Number', how='left')
print(segment)
plt.figure(figsize = (12,4))
sns.scatterplot(x="City", y="Which brand of vehicle do you currently own?", data=segment, hue='Cluster Number',s=400, palette='coolwarm')
plt.title("Simple segment evaluation plot EV", fontsize=15)
plt.xlabel("City", fontsize=10)
plt.ylabel("Which brand of vehicle do you currently own?", fontsize=10)
plt.legend(title='Segment')
plt.show()
   Cluster Number  Do you think Electronic Vehicles are economical?  \
0               0                                          1.720588
1               1                                          1.635220
2               2                                          1.655814
3               3                                          1.634409
4               4                                          1.551948
5               5                                          1.680000

         City  Which brand of vehicle do you currently own?
0  17.838235                                      7.816176
1  18.245283                                      7.301887
2  17.916279                                      8.200000
3  17.806452                                      7.822581
4  16.811688                                      7.740260
5   3.593333                                      7.653333
```
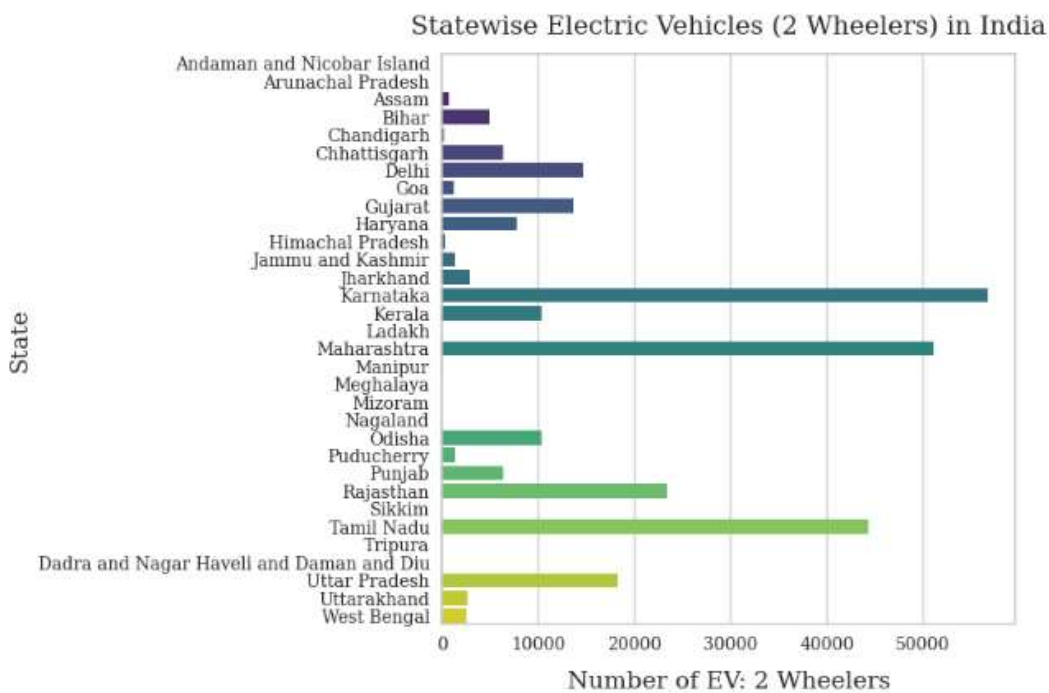


Simple segment evaluation plot EV

**Observations:**

The data points are concentrated in one place. The most people have cars are from Pune city and they most economical cars are Kia and Tata vehicles.

# 2.Dataset 2 -Different types of vehicles and State wise Charging stations



Statewise Electric Vehicles (4 Wheelers) in India



Statewise Electric Vehicles (2 Wheelers) in India
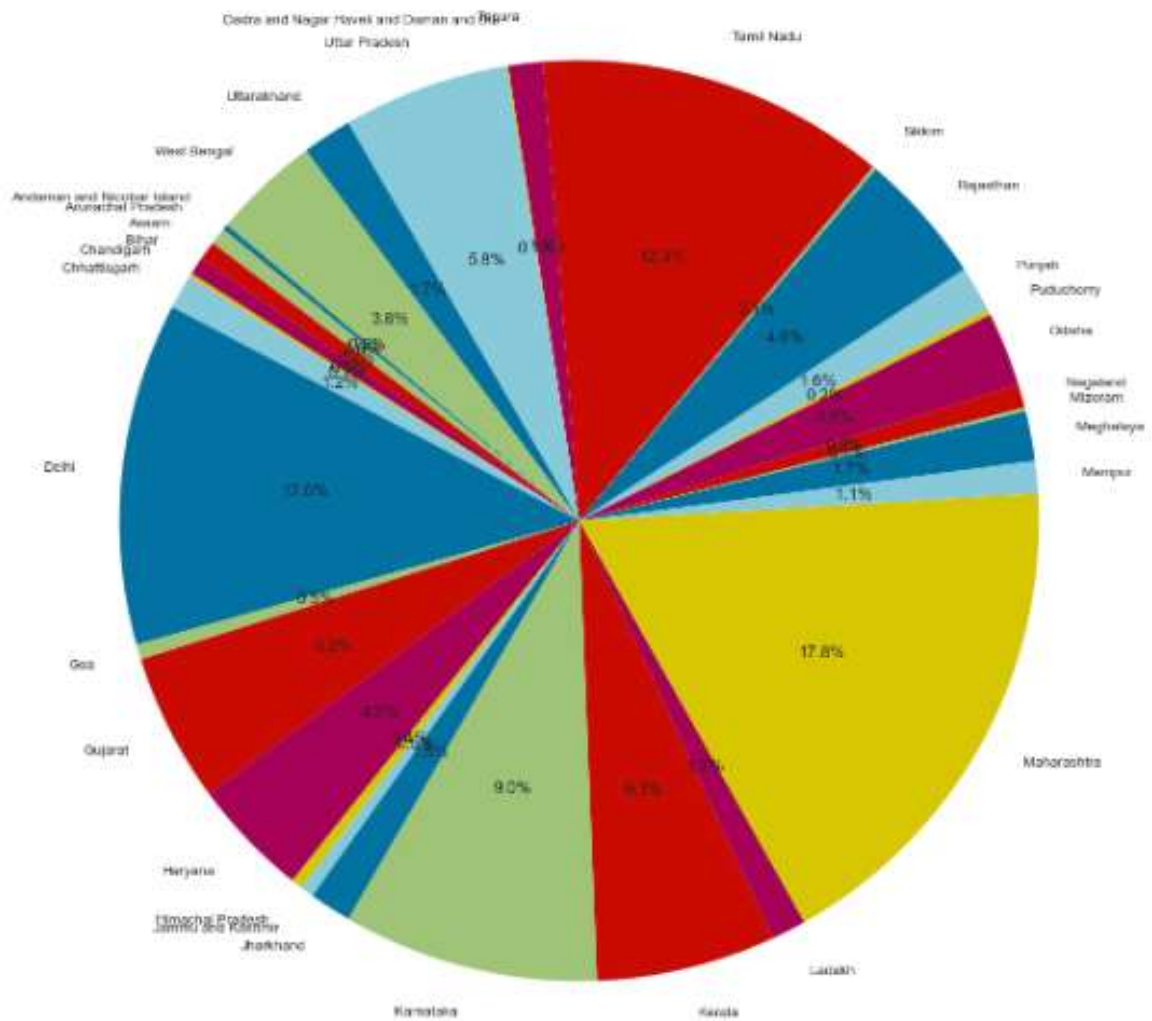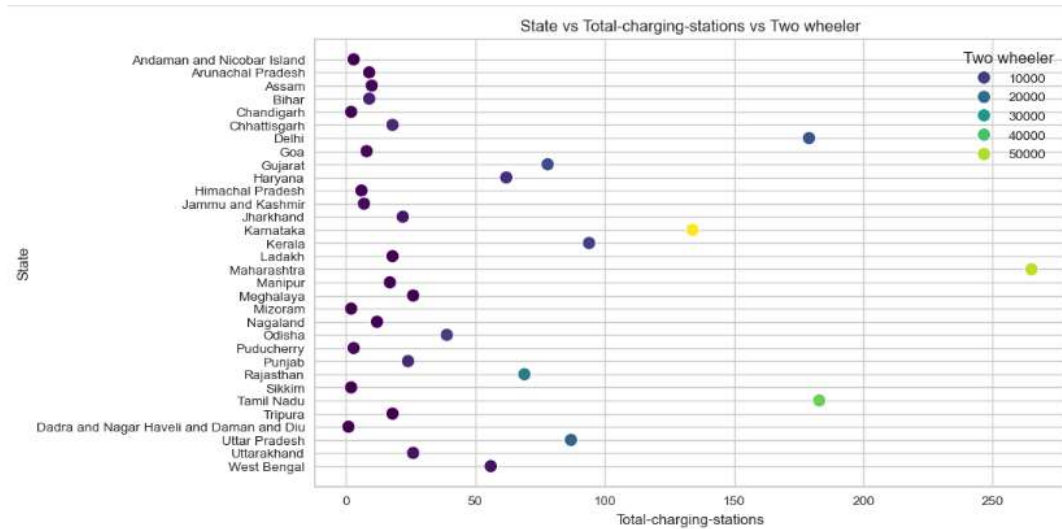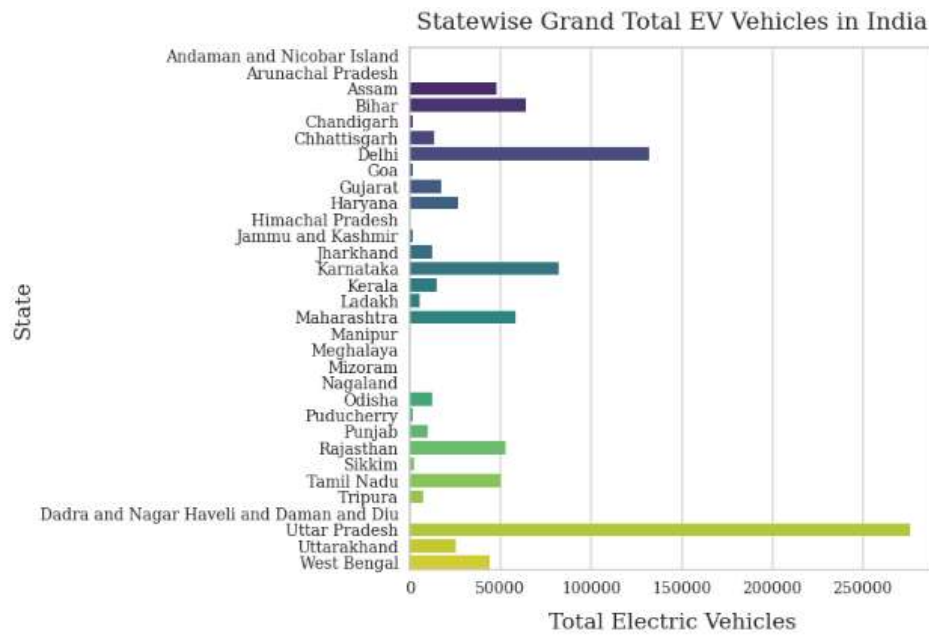
Statewise total-charging-stations in India

Distribution of Charging Stations State-wise

## Statewise Grand Total EV Vehicles in India
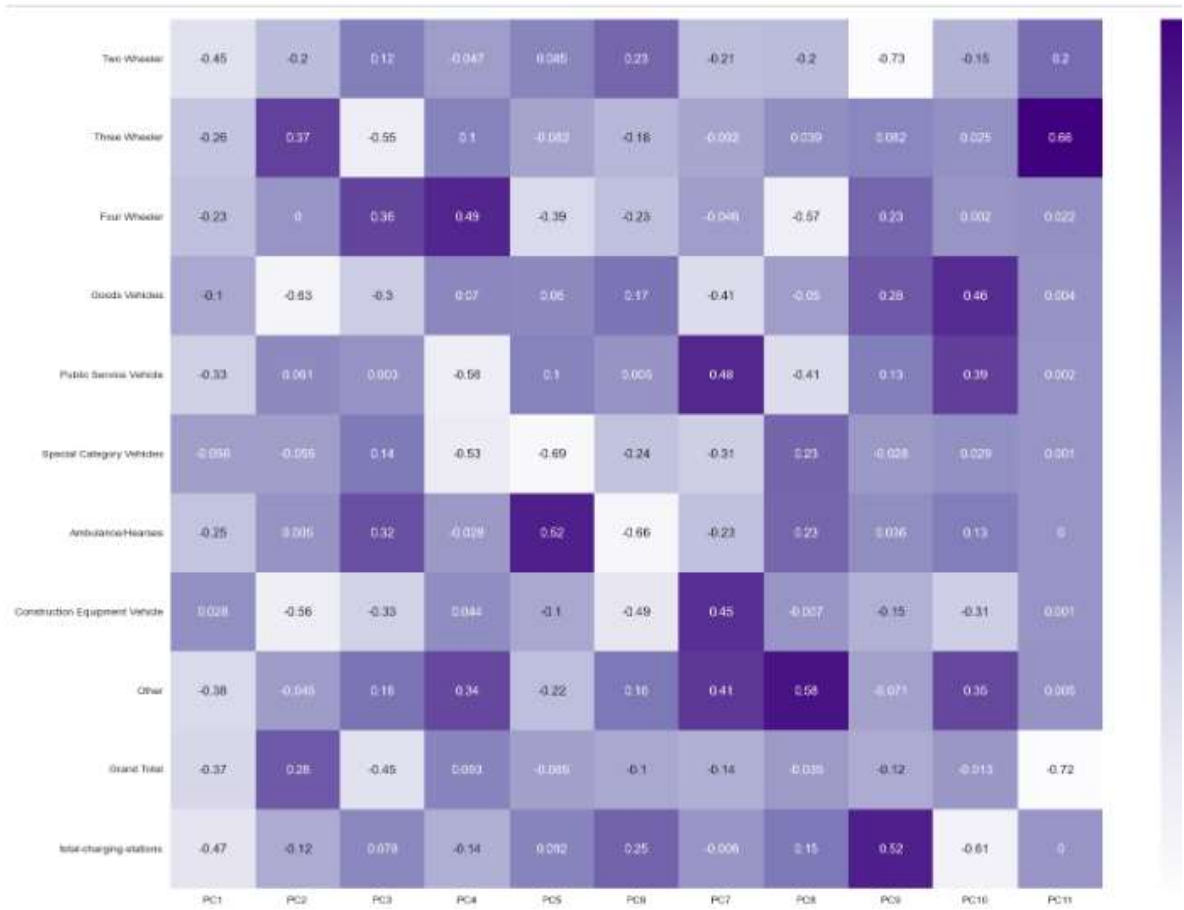


## State vs Total-charging-stations vs Two wheeler



| | PC1 | PC2 | PC3 | PC4 | PC5 | PC6 | PC7 | PC8 | PC9 | PC10 | PC11 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Two Wheeler** | -0.448 | -0.204 | 0.123 | -0.047 | 0.085 | 0.230 | -0.213 | -0.201 | -0.729 | -0.151 | 0.201 |
| **Three Wheeler** | -0.255 | 0.373 | -0.545 | 0.100 | -0.082 | -0.176 | -0.092 | 0.039 | 0.082 | 0.025 | 0.659 |
| **Four Wheeler** | -0.226 | 0.000 | 0.358 | 0.489 | -0.386 | -0.235 | -0.046 | -0.570 | 0.225 | 0.002 | 0.022 |
| **Goods Vehicles** | -0.105 | -0.632 | -0.300 | 0.070 | 0.060 | 0.167 | -0.407 | -0.050 | 0.284 | 0.462 | 0.004 |
| **Public Service Vehicle** | -0.325 | 0.061 | 0.003 | -0.558 | 0.105 | 0.005 | 0.484 | -0.410 | 0.126 | 0.387 | 0.002 |
| **Special Category Vehicles** | -0.056 | -0.055 | 0.140 | -0.534 | -0.691 | -0.242 | -0.312 | 0.235 | -0.028 | 0.029 | 0.001 |
| **Ambulance/Hearses** | -0.248 | 0.005 | 0.317 | -0.028 | 0.523 | -0.663 | -0.232 | 0.226 | 0.036 | 0.131 | 0.000 |
| **Construction Equipment Vehicle** | 0.028 | -0.565 | -0.331 | 0.044 | -0.102 | -0.486 | 0.452 | -0.007 | -0.148 | -0.310 | 0.001 |
| **Other** | -0.382 | -0.045 | 0.176 | 0.340 | -0.219 | 0.158 | 0.411 | 0.577 | -0.071 | 0.354 | 0.005 |
| **Grand Total** | -0.368 | 0.278 | -0.451 | 0.093 | -0.065 | -0.103 | -0.141 | -0.035 | -0.120 | -0.013 | -0.725 |
| **total-charging-stations** | -0.469 | -0.120 | 0.078 | -0.141 | 0.092 | 0.251 | -0.006 | 0.150 | 0.521 | -0.612 | 0.000 |

| | PC1 | PC2 | PC3 | PC4 | PC5 | PC6 | PC7 | PC8 | PC9 | PC10 | PC11 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Two Wheeler | -0.45 | -0.2 | 0.12 | -0.047 | 0.085 | 0.23 | -0.21 | -0.2 | -0.73 | -0.15 | 0.2 |
| Three Wheeler | -0.26 | 0.37 | -0.55 | 0.1 | -0.082 | -0.16 | -0.002 | 0.039 | 0.082 | 0.025 | 0.66 |
| Four Wheeler | -0.23 | 0 | 0.36 | 0.49 | -0.39 | -0.23 | -0.046 | -0.57 | 0.23 | 0.002 | 0.022 |
| Goods Vehicles | -0.1 | -0.63 | -0.3 | 0.07 | 0.06 | 0.17 | -0.41 | -0.05 | 0.28 | 0.46 | 0.004 |
| Public Service Vehicle | -0.33 | 0.061 | 0.003 | -0.56 | 0.1 | 0.005 | 0.48 | -0.41 | 0.13 | 0.39 | 0.002 |
| Special Category Vehicles | -0.056 | -0.055 | 0.14 | -0.53 | -0.69 | -0.24 | -0.31 | 0.23 | -0.028 | 0.029 | 0.001 |
| Ambulance/Hearses | -0.25 | 0.005 | 0.32 | -0.028 | 0.52 | -0.66 | -0.23 | 0.23 | 0.036 | 0.13 | 0 |
| Construction Equipment Vehicle | 0.028 | -0.56 | -0.33 | 0.044 | -0.1 | -0.49 | 0.45 | -0.007 | -0.15 | -0.31 | 0.001 |
| Other | -0.38 | -0.046 | 0.18 | 0.34 | -0.22 | 0.16 | 0.41 | 0.58 | -0.071 | 0.35 | 0.005 |
| Grand Total | -0.37 | 0.28 | -0.45 | 0.093 | -0.085 | -0.1 | -0.14 | 0.039 | -0.12 | -0.013 | -0.72 |
| total-charging-stations | -0.47 | -0.12 | 0.079 | -0.14 | 0.092 | 0.25 | -0.006 | 0.15 | 0.52 | -0.61 | 0 |



Distortion Score Elbow for KMeans Clustering

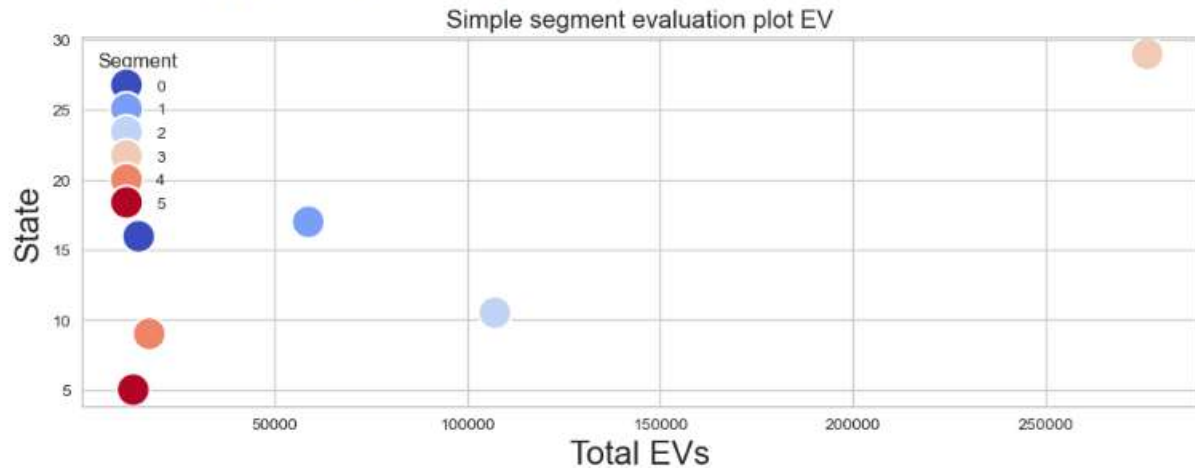--- elbow at k = 6, score = 98.794

[136]: <Axes: title={'center': 'Distortion Score Elbow for KMeans Clustering'}, xlabel='k', ylabel='distortion score'>

```
segment = four_wheeler.merge(total, on='Cluster Number', how='left').merge(state, on='Cluster Number', how='left')
print(segment)
plt.figure(figsize = (12,4))
sns.scatterplot(x="Grand Total", y="State Name", data=segment, hue='Cluster Number',s=400, palette='coolwarm')
plt.title("Simple segment evaluation plot EV", fontsize=15)
plt.xlabel("Total EVs", fontsize=20)
plt.ylabel("State", fontsize=20)
plt.legend(title='Segment')
plt.show()
```

```
   Cluster Number  Four Wheeler   Grand Total  State Name
0               0    549.076923  14844.692308   15.961538
1               1      2.000000  58815.000000   17.000000
2               2   5131.500000 107174.000000   10.500000
3               3    368.000000 276217.000000   29.000000
4               4   1309.000000  17593.000000    9.000000
5               5    117.000000  13428.000000    5.000000
```
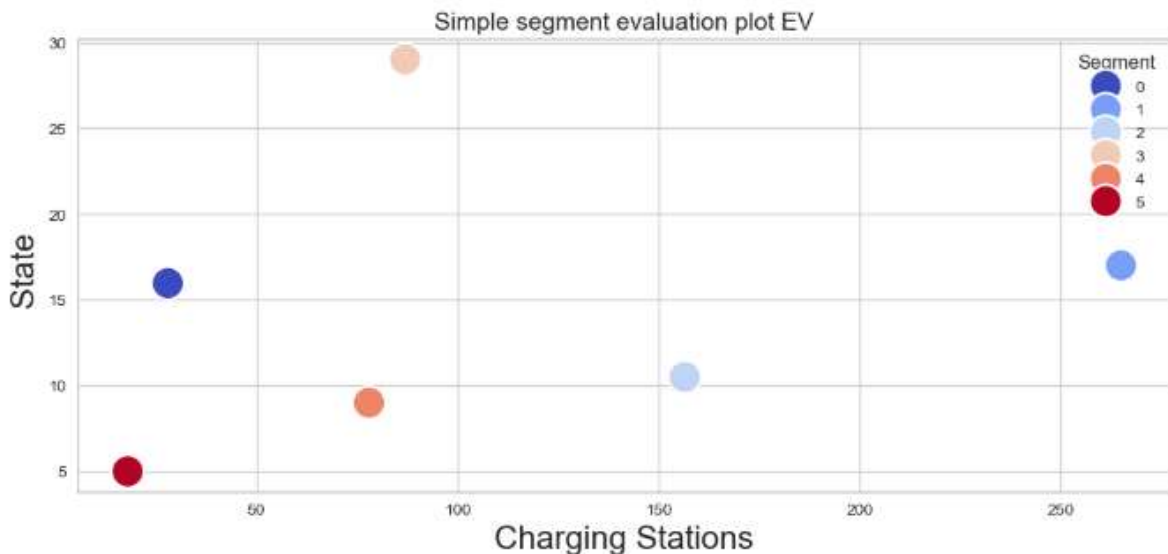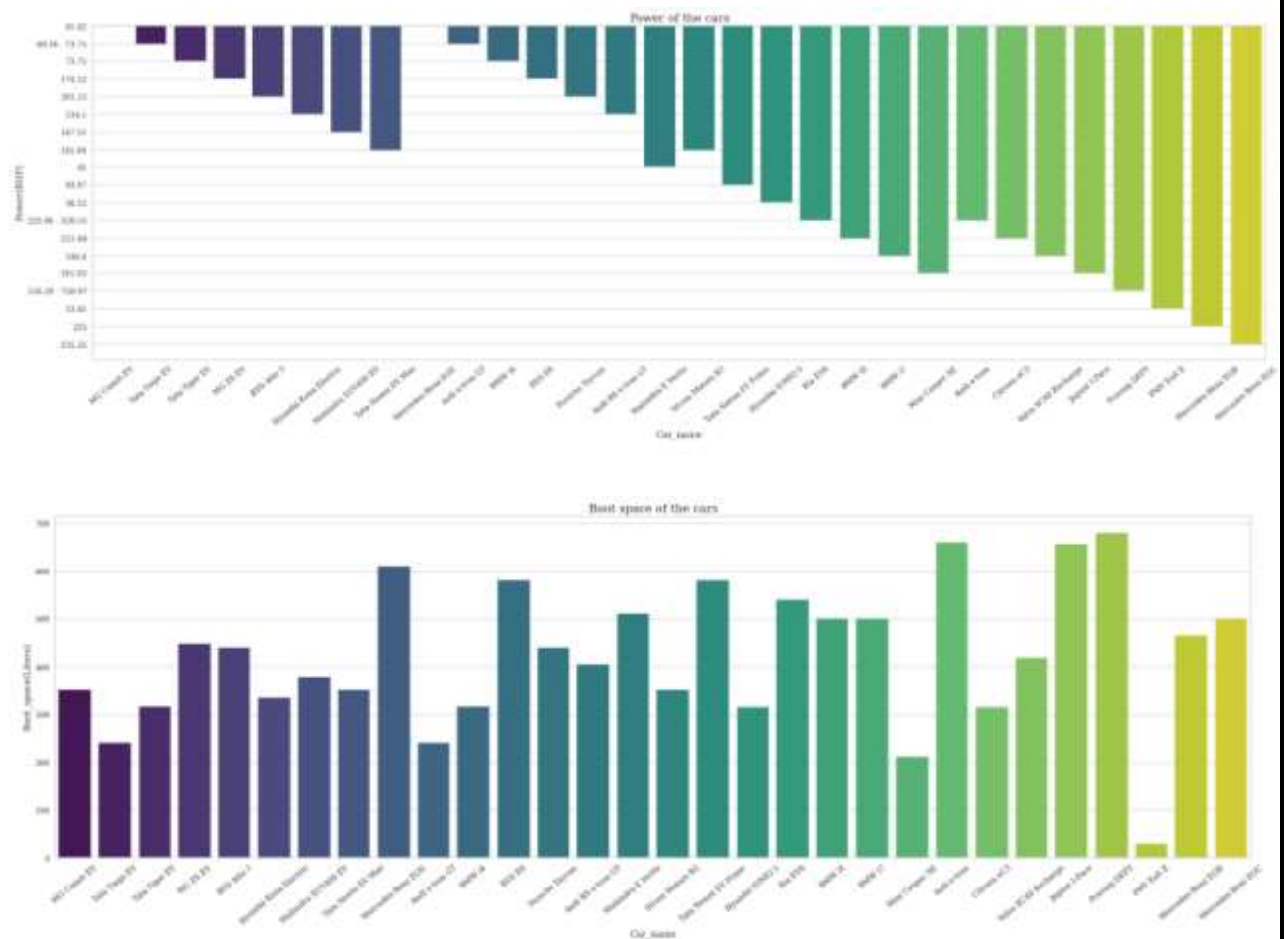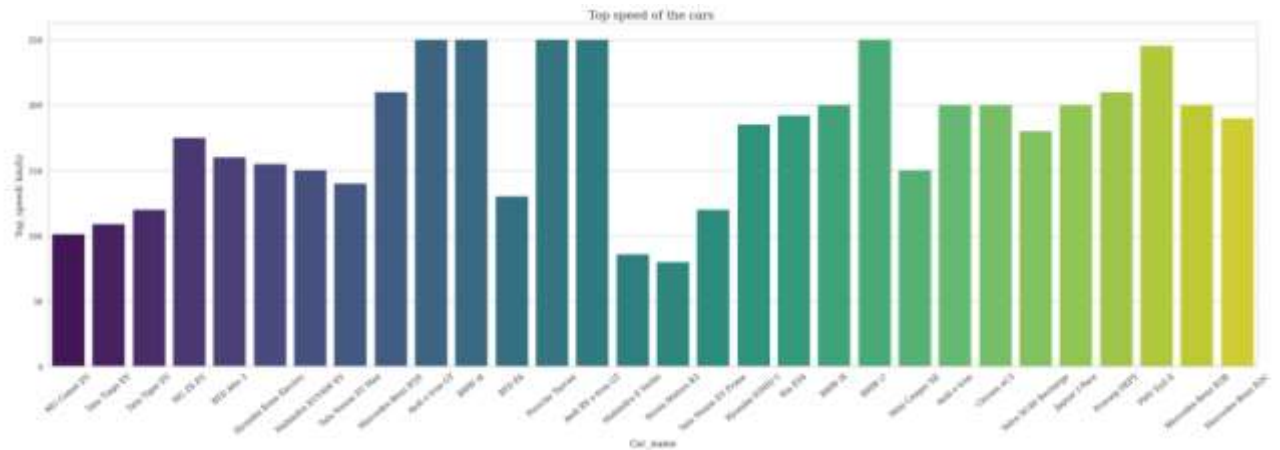


Simple segment evaluation plot EV

```
segment = total.merge(charging_stations, on='Cluster Number', how='left').merge(state, on='Cluster Number', how='left')
print(segment)
plt.figure(figsize = (12,5))
sns.scatterplot(x="total-charging-stations", y="State Name", data=segment, hue='Cluster Number',s=400, palette='coolwarm')
plt.title("Simple segment evaluation plot EV", fontsize=15)
plt.xlabel("Charging Stations", fontsize=20)
plt.ylabel("State ", fontsize=20)
plt.legend(title='Segment')
plt.show()
```

```
   Cluster Number   Grand Total  total-charging-stations  State Name
0               0  14844.692308                     28.0   15.961538
1               1  58815.000000                    265.0   17.000000
2               2 107174.000000                    156.5   10.500000
3               3 276217.000000                     87.0   29.000000
4               4  17593.000000                     78.0    9.000000
5               5  13428.000000                     18.0    5.000000
```



Simple segment evaluation plot EV

**Observations:**

Karnataka, Ladakh have the maximum number of 4-wheeler EVs.

Karnataka, Maharashtra and Tamil Nadu have the maximum number of 2-wheeler EVs.

Maharashtra, Tamil Nadu, Delhi have the maximum number of charging stations.

Maharashtra has 17.8% of charging stations of India.

Total number of EVs are more in UP.

# 3.Dataset 3- EV cars of India

Top speed of the cars

**Observations:**

Mercedes-Volvo, Pravaig DEFY, Benz, BYD, Audi tops the list of EVs with the maximum power in the Indian automobile market.

Pravaig DEFY, Volvo and Audi have more boot space.

Mercedes-Benz, Audi, BMW, Porsche Taycan have maximum speed.

## Conclusion

Based on the market segmentation analysis, we identified distinct customer segments with unique characteristics and preferences. These segments can be targeted with tailored marketing strategies to enhance customer engagement and drive sales. By understanding these segments, businesses can optimize their product offerings, pricing strategies, and marketing campaigns to better meet the needs of their target audience. This approach not only improves customer satisfaction but also increases market share and profitability.

1. By using the above 3 datasets, we can conclude that Maharashtra is the best State to open new Electric Vehicle Startup in India with charging maximum number of charging stations. In Maharashtra, Pune is the best City to open the start up. After Maharashtra we can select Karnataka, Delhi, Tamil Nadu.

2. We can choose 2-wheeler and 4-wheeler market as preferred EVs from the data. In 4-wheeler category, SUVs are most preferable type.

3. For top model EV vehicles Pravaig DEFY, Benz, BYD, Audi are the best. For medium priced vehicles Tata is the best option.