

## Data pre-processing

### Required libraries

In order to perform EDA and clustering on the collected data, the following Python libraries are used:

1. Pandas: for data handling/manipulation
2. Matplotlib and Seaborn: for data visualization
3. Scikit-learn: for the k-means clustering algorithm and some other algorithms

```
|: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import seaborn as sb
import statsmodels.api as sm
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
import plotly.io as pio
pio.renderers.default = "svg"
from sklearn.cluster import KMeans
```

### Pulling the datasets

#### [Dataset 1](#)

```
df1 = pd.read_csv('charging station.csv')
```

```
df1.head()
```

	Sl. no.	State/UT	City	Operational Charging stations Under FAME-II
0	1	Delhi	Delhi	15
1	2	Maharashtra	Navi Mumbai	1
2	3	Maharashtra	Nagpur	7
3	4	Tamil Nadu	Chennai	8
4	5	Kerala	Thrissur	8

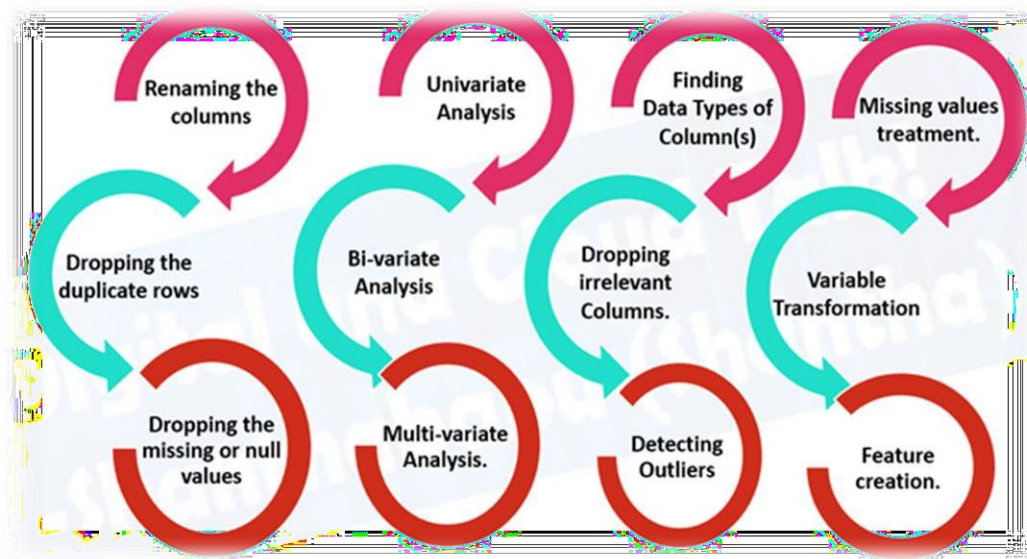
## Dataset 2

```
df2 = pd.read_csv('data.csv')
df2.drop('Unnamed: 0', axis=1, inplace=True)
df2['inr(10e3)'] = df2['PriceEuro']*0.08320
df2['RapidCharge'].replace(to_replace=['No', 'Yes'], value=[0, 1], inplace=True)
df2.head()
```

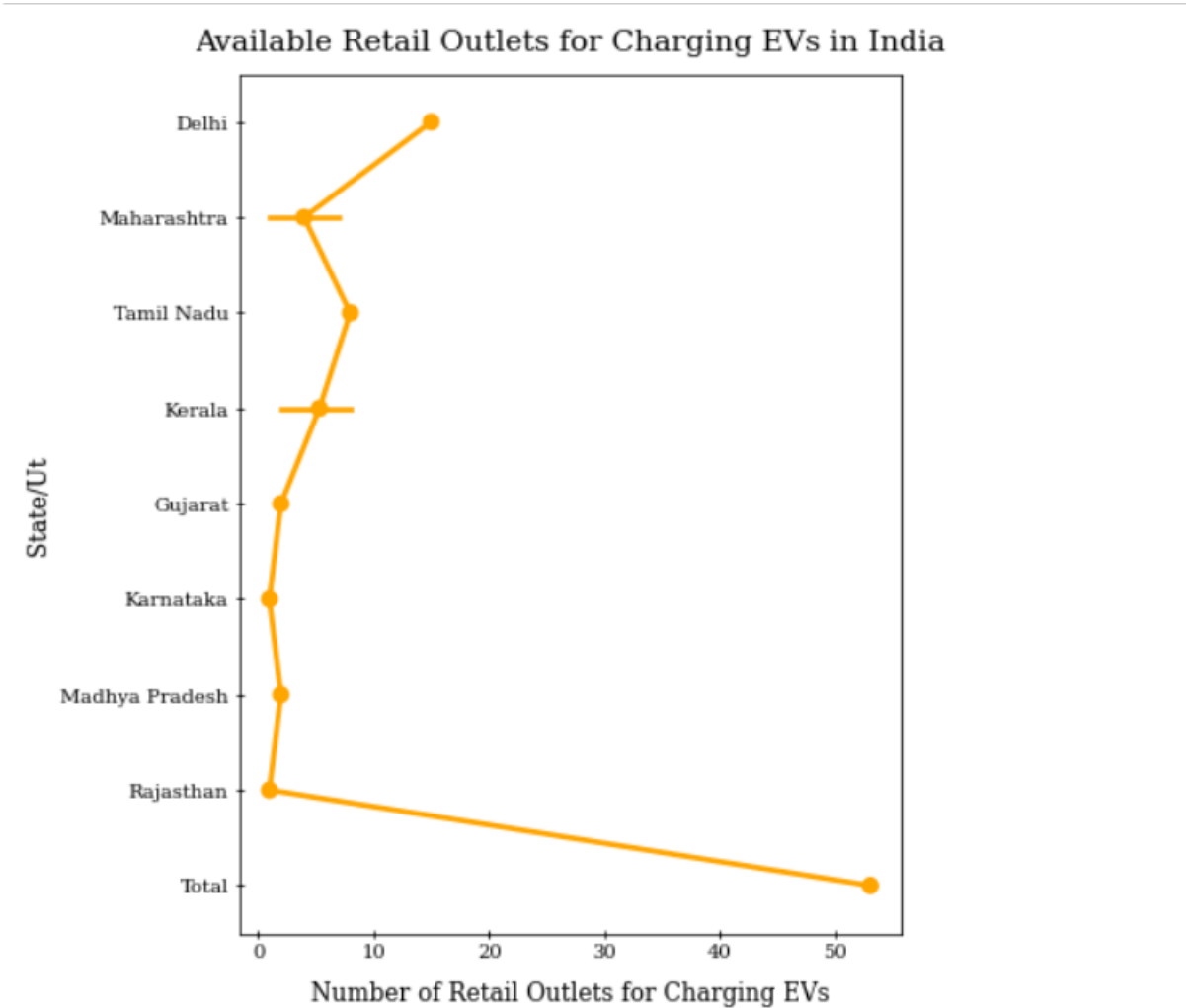
	Brand	Model	AccelSec	TopSpeed_KmH	Range_Km	Efficiency_WhKm	FastCharge_KmH	RapidCharge	PowerTrain	PlugType	BodyStyle	Segment	Score
0	Tesla	Model 3 Long Range Dual Motor	4.6	233	450	161	940	1	AWD	Type 2 CCS	Sedan	D	8.5
1	Volkswagen	ID.3 Pure	10.0	160	270	167	250	0	RWD	Type 2 CCS	Hatchback	C	7.5
2	Polestar	2	4.7	210	400	181	620	1	AWD	Type 2 CCS	Liftback	D	8.0
3	BMW	iX3	6.8	180	360	206	560	1	RWD	Type 2 CCS	SUV	D	8.0
4	Honda	e	9.5	145	170	168	190	1	RWD	Type 2 CCS	Hatchback	B	6.5

## Exploratory Data Analysis

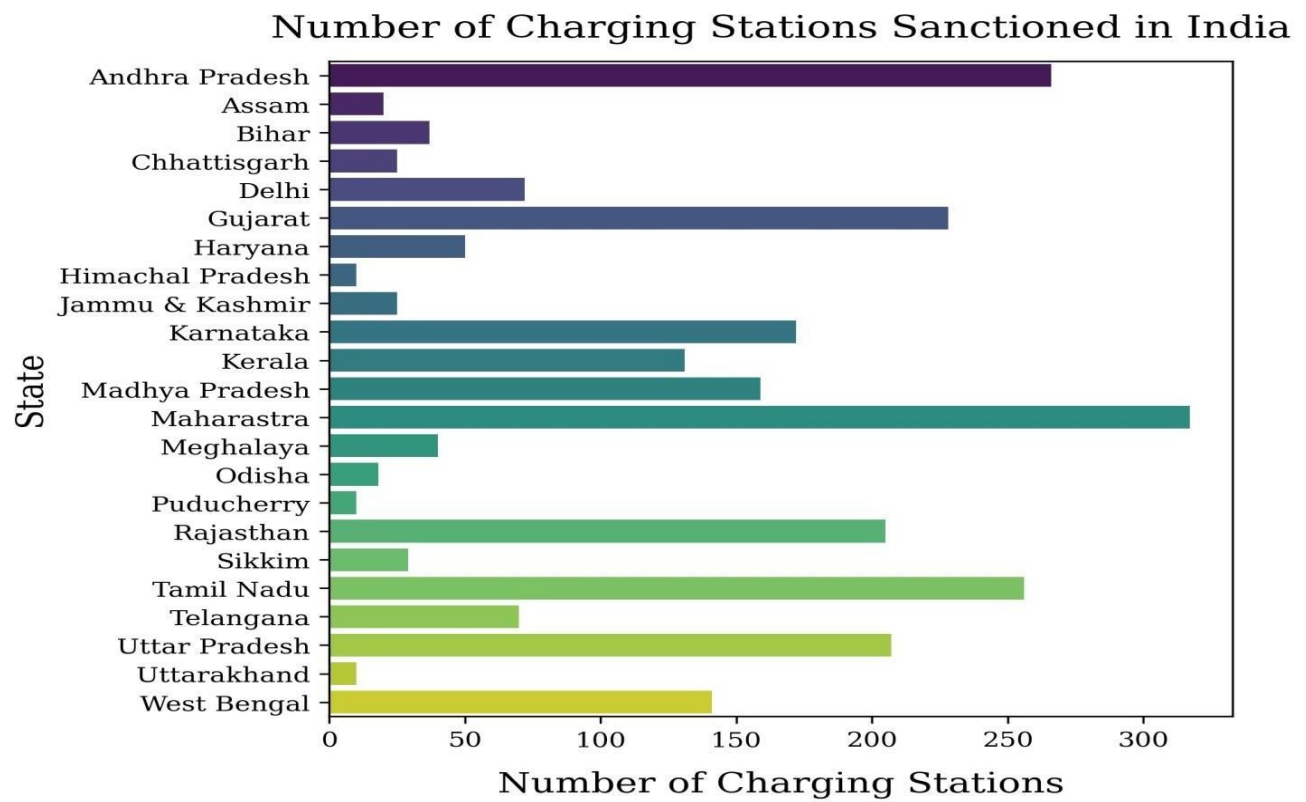
Exploratory Data Analysis, popularly abbreviated as EDA, is one of the most important steps in the data science pipeline. It is the process of gaining the information present inside the data with the help of summary statistics and visual representations. Keys features of this technique are presented in the below image.



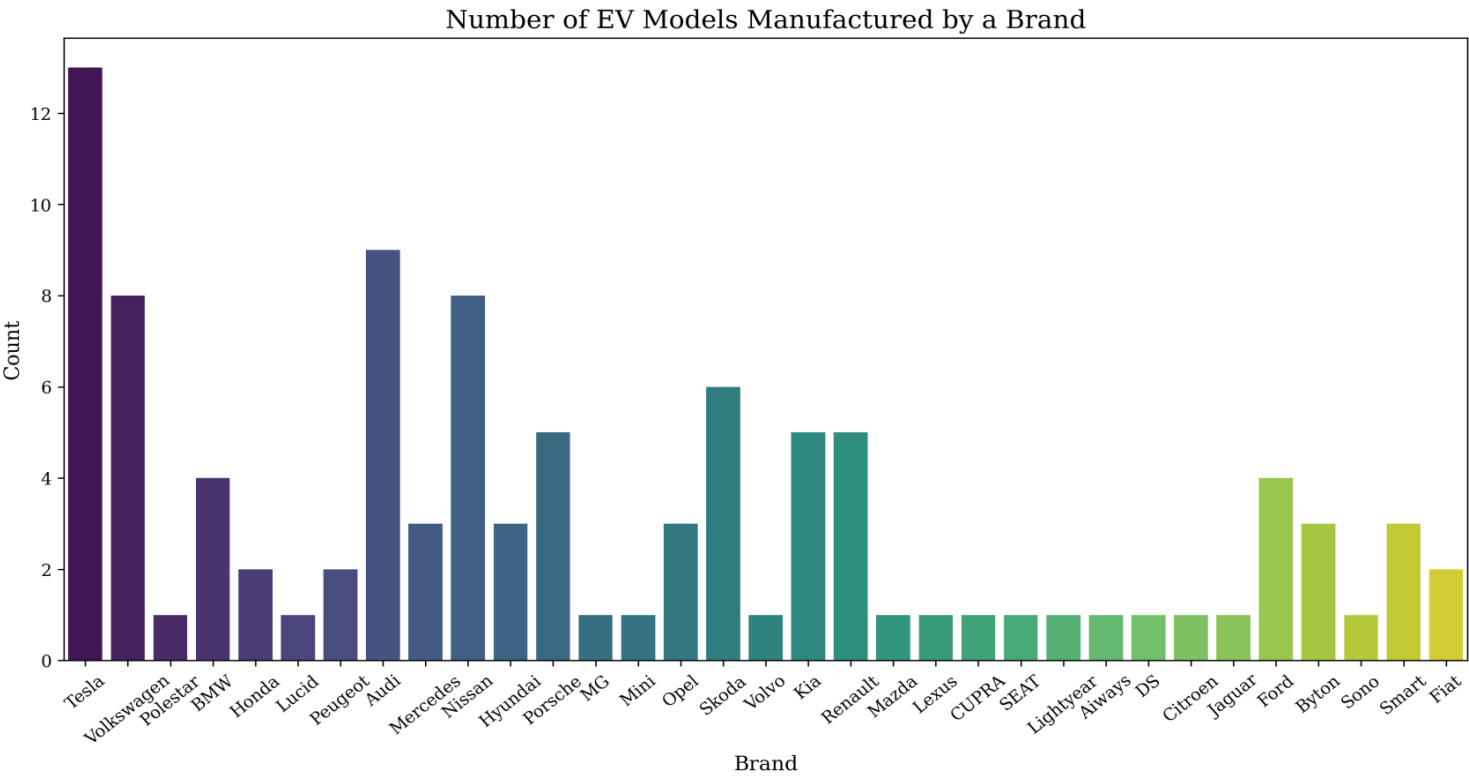
*Retail outlets in India for charging EVs*



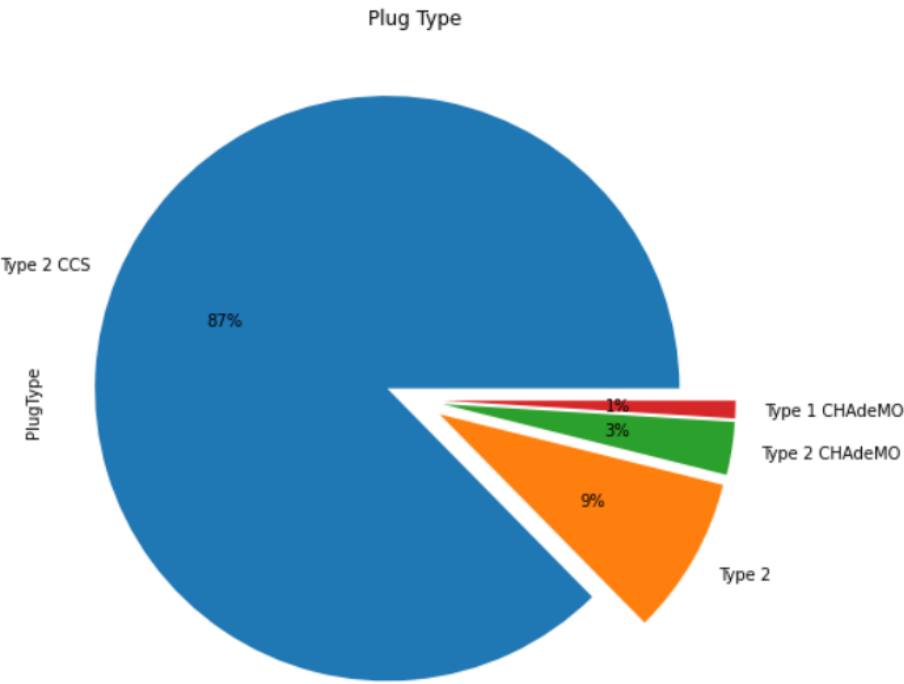
*Number of charging stations sanctioned by Government of India*



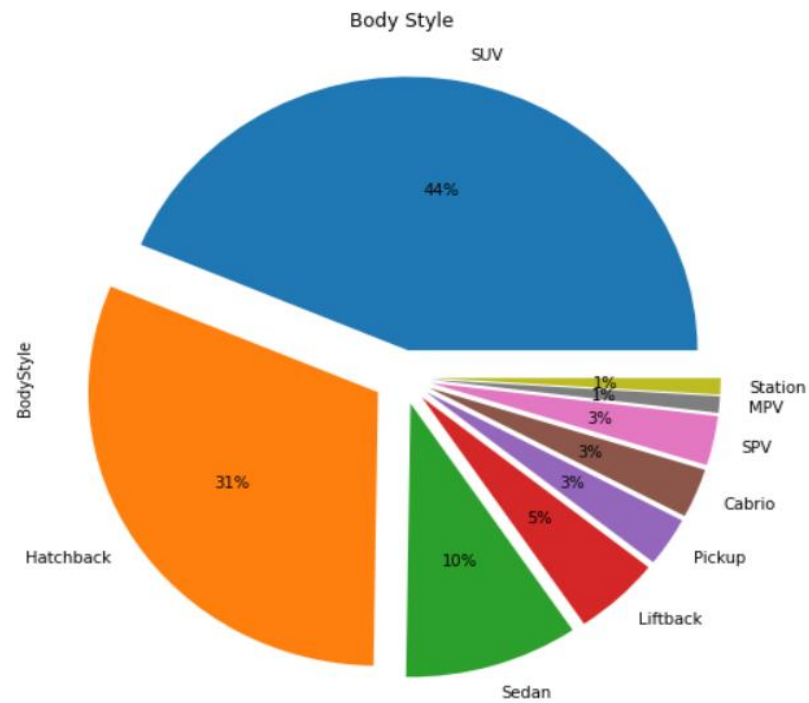
Top EV manufacturing brands in India



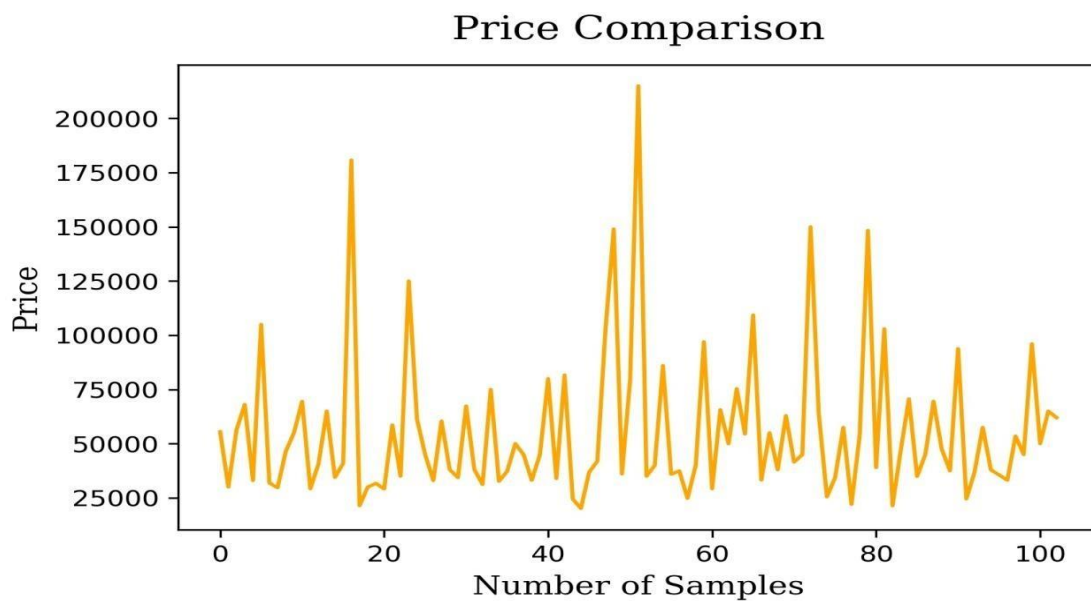
Types of EV plugs available in India



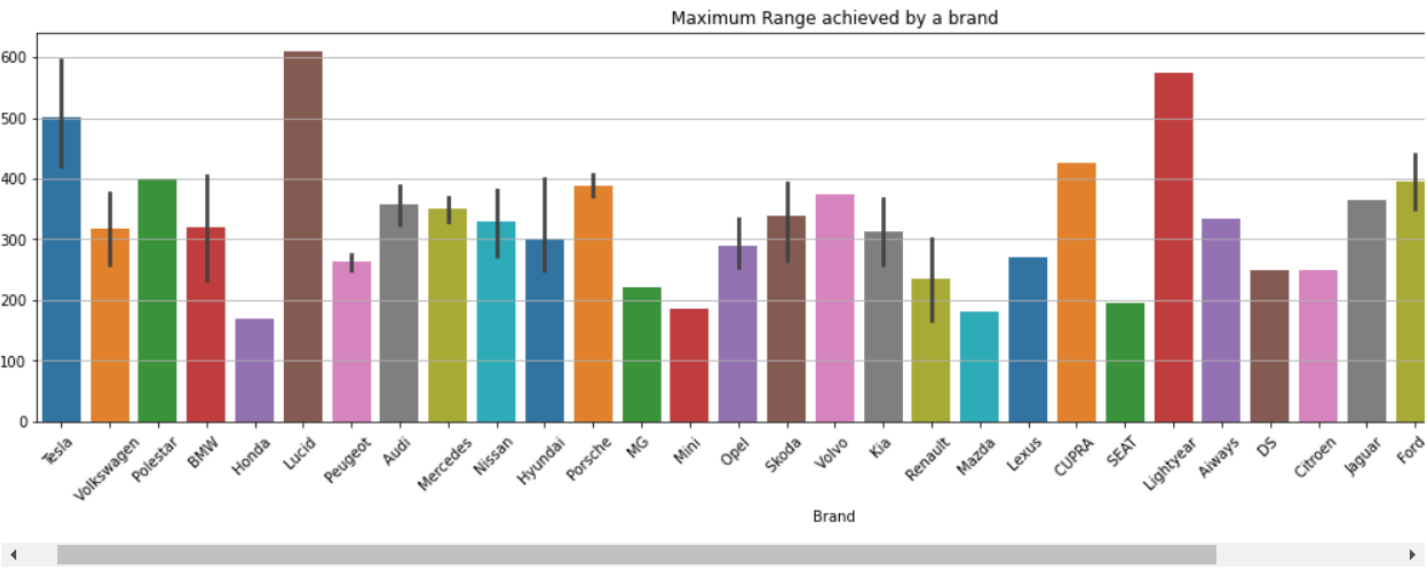
### *Body types of EVs in India*



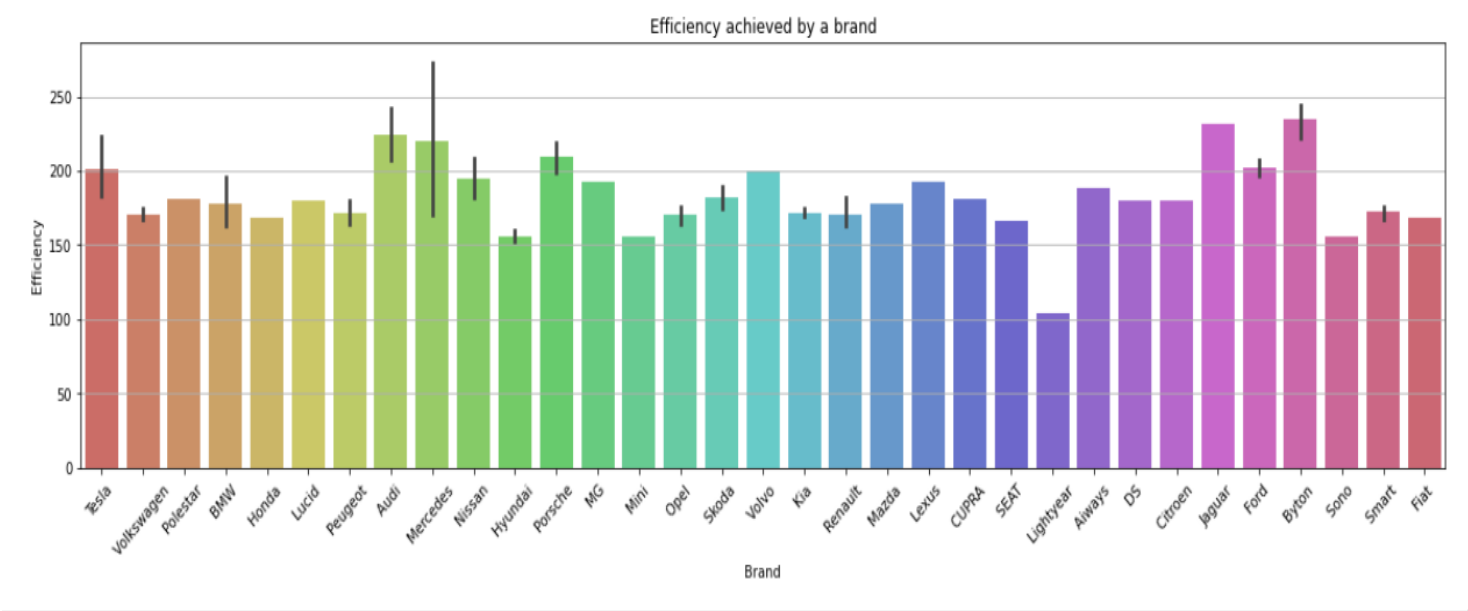
### *Price comparison of different brands of EVs in India*



Efficiency achieved by a brand

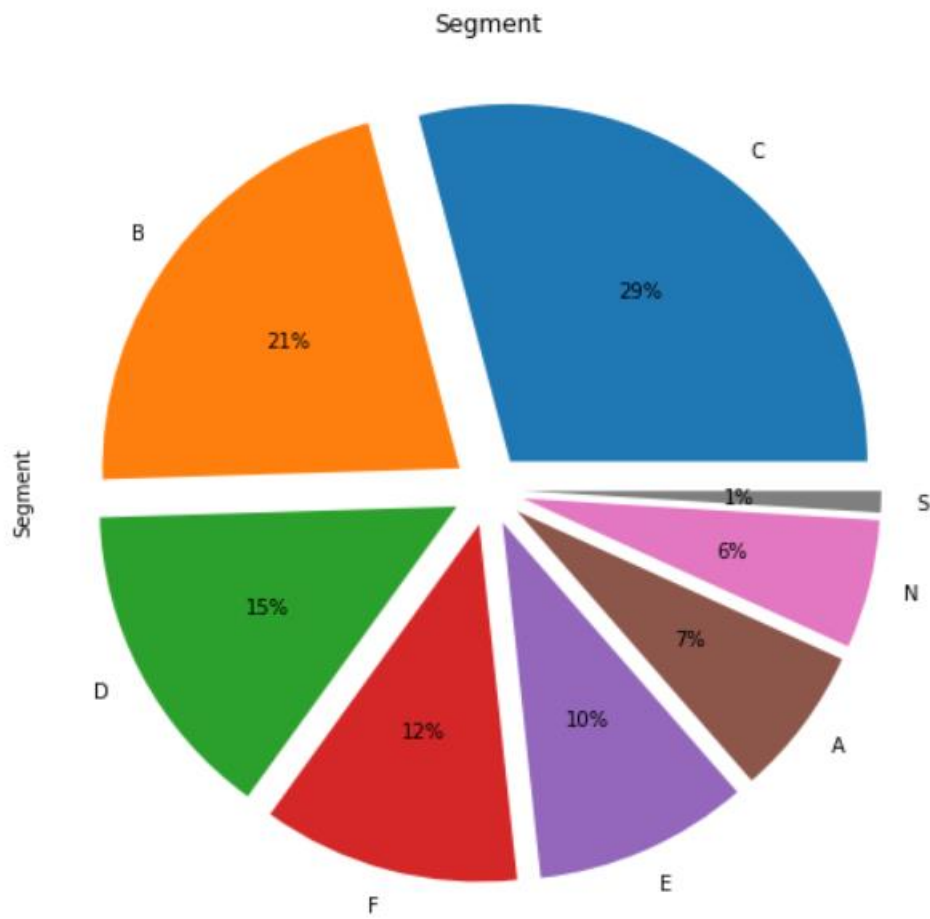


Maximum Range achieved by a brand





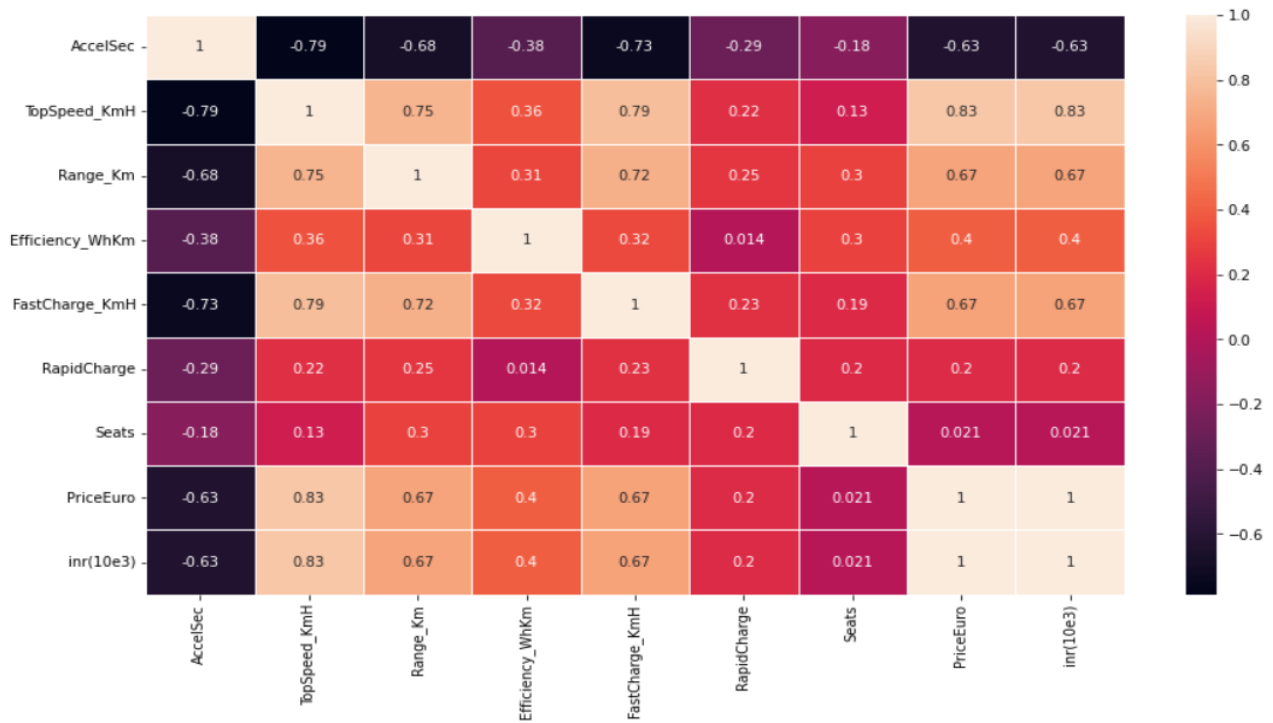
*EV Segments in India*



## Correlation Matrix

```
: # Heatmap to show the correlation of the dataset 2
ax= plt.figure(figsize=(15,8))
sb.heatmap(df2.corr(),linewidths=1,linecolor='white',annot=True)

: <AxesSubplot:>
```



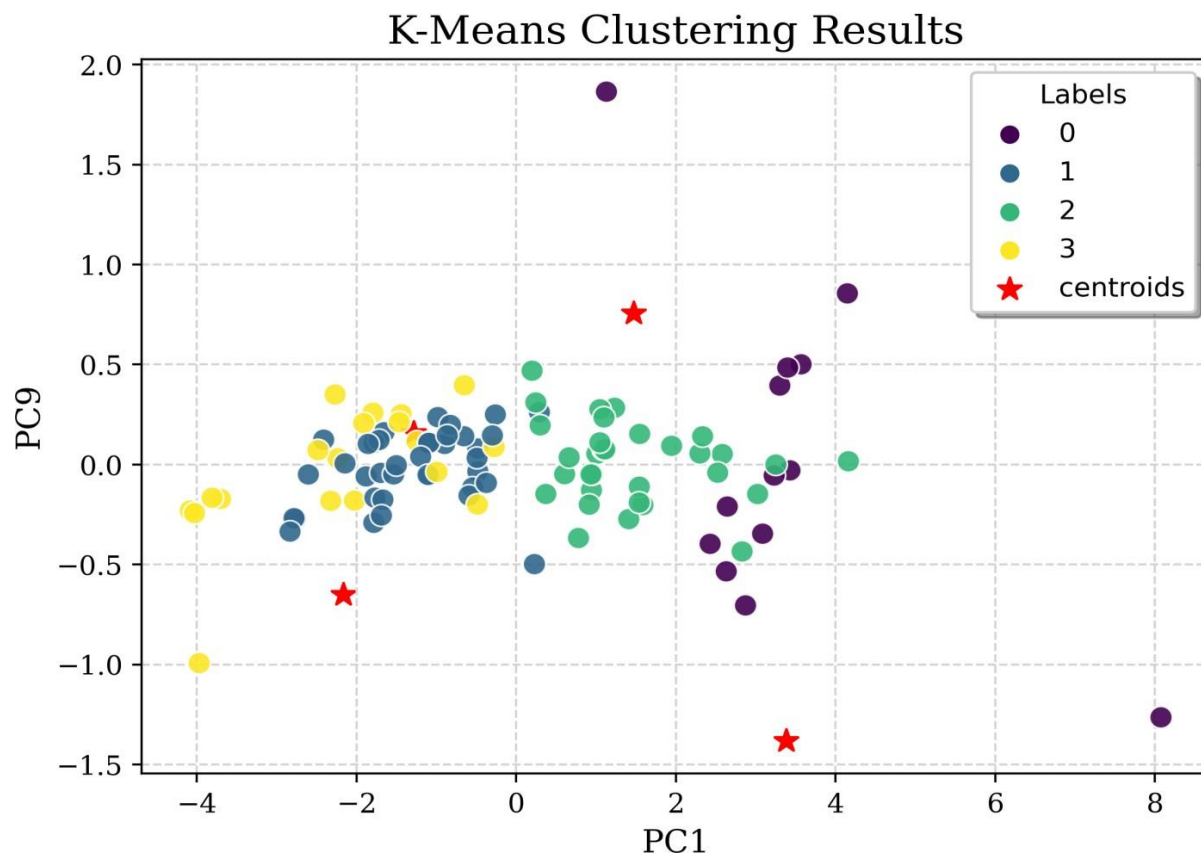
## Segmentation Approaches

### Clustering

Clustering is an unsupervised machine learning technique of grouping similar data points into clusters. The sole objective of this technique is to segregate datapoints with similar traits and place them into different clusters. There are several algorithms to perform clustering on data such as k-means clustering, hierarchical clustering, density-based clustering etc.

### K-Means Clustering

K-Means Clustering is an unsupervised learning algorithm whose job is to group the unlabelled dataset into different clusters where each datapoint belongs to only one cluster. Here, K is the number of clusters that need to be created in the process. The algorithm finds its applicability into a variety of use cases including market segmentation, image segmentation, image compression, document clustering etc. The below image is the results of clustering on one of our datasets.



The K-Means Algorithm works the following way:

1. Specify the number of clusters, i.e. K
2. Select K random points in the dataset. These points will be the centroids (centres) of each of the K clusters.
3. Assign each data point in the dataset to one of the K centroids, based on its distance from each of the centroids.
4. Consider this clustering to be correct and reassign the Centroids to the mean of these clusters.
5. Repeat Step 3. If any of the points change clusters, Go to step 4. Else Go to step 6.
6. Calculate the variance of each of the clusters.
7. Repeat this clustering 'n' number of times until the sum of variance of each cluster is minimum.

### Principle Component Analysis

Principal component analysis (PCA) is a linear dimensionality-reduction technique that is used to reduce the dimensionality of large data sets by transforming a large set of variables into a smaller one while preserving most of the information present in the large set.

### Elbow Method

The Elbow method is a way of determining the optimal number of clusters (k) in K-Means Clustering. It is based on calculating the Within Cluster Sum of Squared Errors (WCSS) for a different number of clusters (k) and selecting the k for which change in WCSS first starts to diminish. When you plot its graph, at one point the line starts to run parallel to the X-axis and that point, known as the Elbow Point, is considered as the best value for the k (as 4 in the below figure).

