

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

```
[1] # importing the dependencies
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans
```

Pulling the datasets

[Dataset 1](#)

```
✓ [2] # fetching dataset - 1
0s df1 = pd.read_csv('/content/ev_charger.csv')
df1.head()
```

	Region	2W	3W	4W	Bus	Chargers
0	Uttar Pradesh	9852	42881	458	197	207
1	Maharastra	38558	893	1895	186	317
2	Karnataka	32844	568	589	57	172
3	Tamil Nadu	25642	396	426	0	256
4	Gujarat	22359	254	423	22	228

Dataset 2

```
✓ [4] # fetching dataset - 2
0s df2 = pd.read_excel('/content/Ev charging station data.xlsx', sheet_name='Table 4', header=1)
df2.head()
```

	State/UT	EV Charging Facility
0	Andhra Pradesh	65
1	Arunachal Pradesh	4
2	Assam	19
3	Bihar	26
4	Chandigarh	4

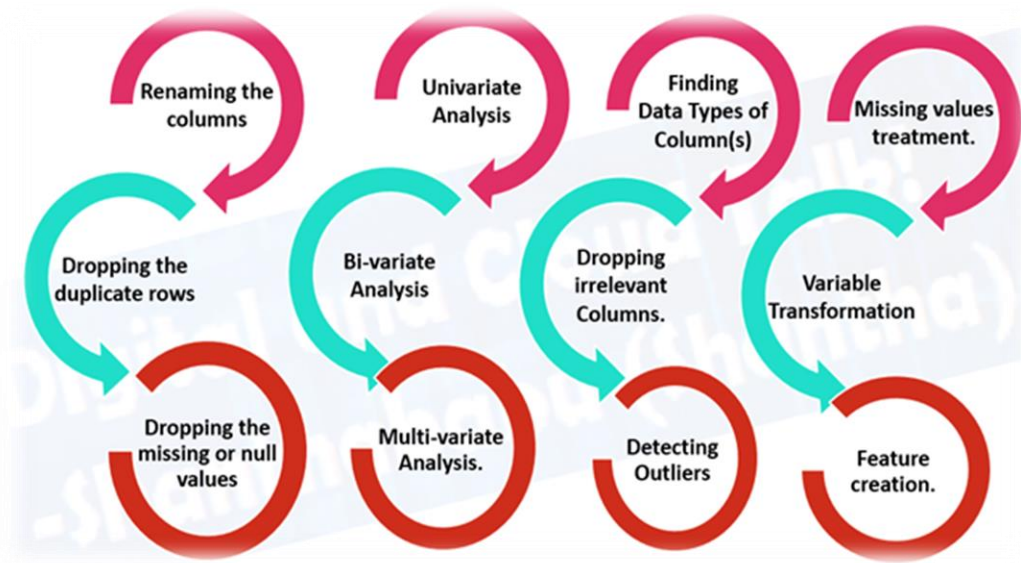
Dataset 3

```
✓ [5] # fetching dataset - 3
1s df3 = pd.read_excel('/content/ev_market_india_data.xlsx')
df3.head()
```

	Brand	Model	AccelSec	TopSpeed_KmH	Range_Km	Efficiency_WhKm	FastCharge_KmH	RapidCharge	PowerTrain	PlugType	BodyStyle	Segment	Seats	PriceEuro
0	Tesla	Model 3 Long Range Dual Motor	4.6	233	450	161	940	Yes	AWD	Type 2 CCS	Sedan	D	5	55480
1	Volkswagen	ID.3 Pure	10.0	160	270	167	250	No	RWD	Type 2 CCS	Hatchback	C	5	30000
2	Polestar	2	4.7	210	400	181	620	Yes	AWD	Type 2 CCS	Liftback	D	5	56440
3	BMW	ix3	6.8	180	360	206	560	Yes	RWD	Type 2 CCS	SUV	D	5	68040
4	Honda	e	9.5	145	170	168	190	Yes	RWD	Type 2 CCS	Hatchback	B	4	32997

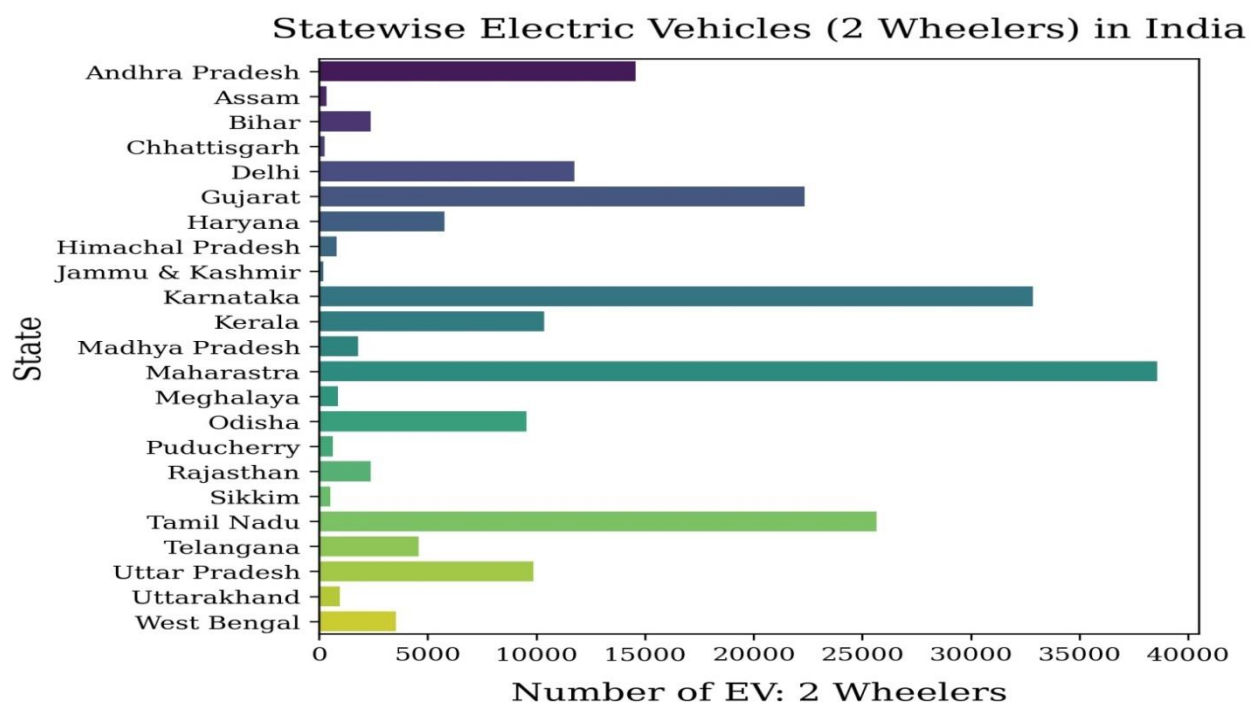
Exploratory Data Analysis

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.

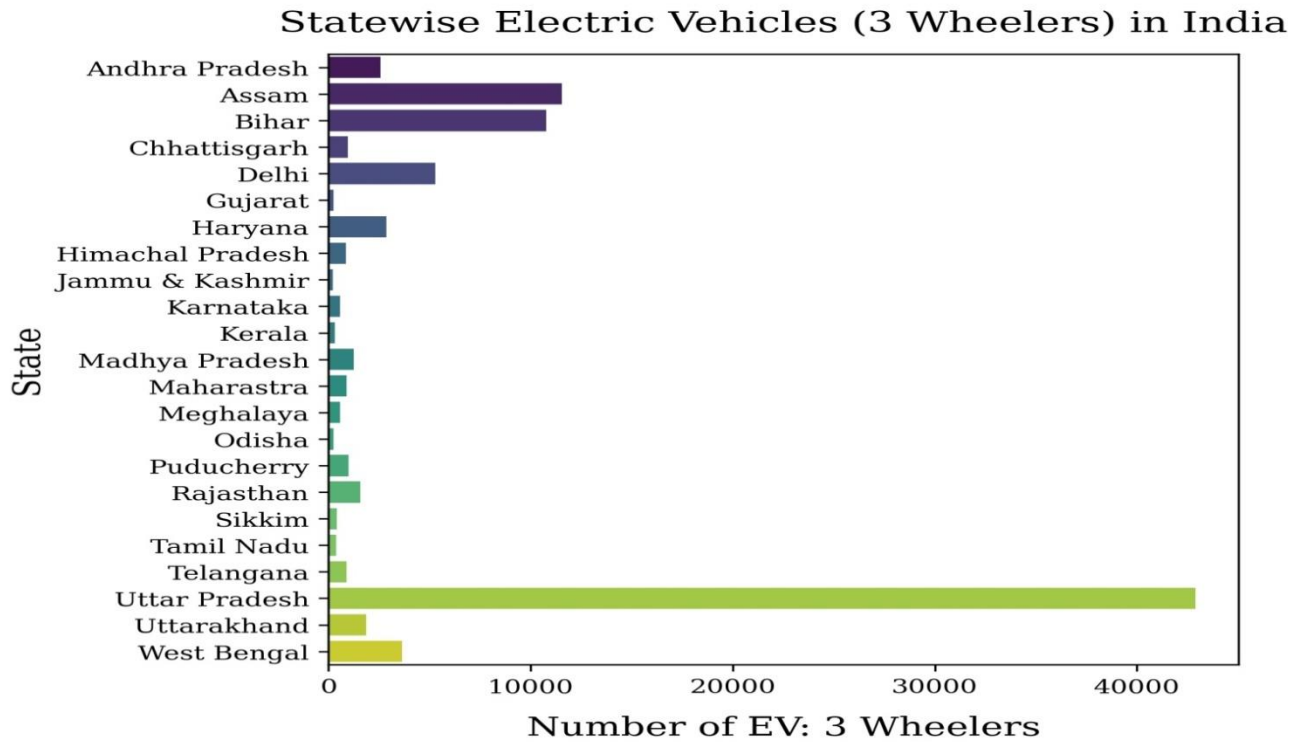


Implementing EDA on the datasets

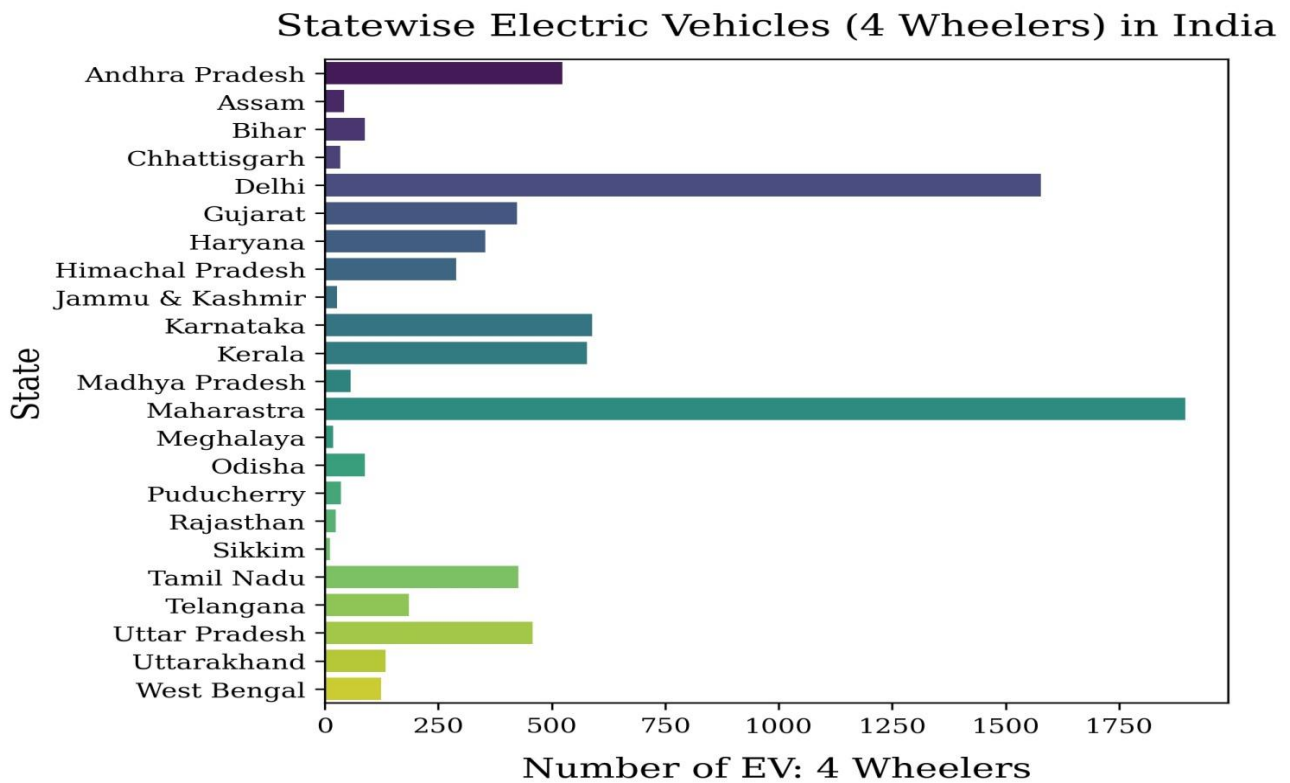
Number of 2-wheeler EVs in India



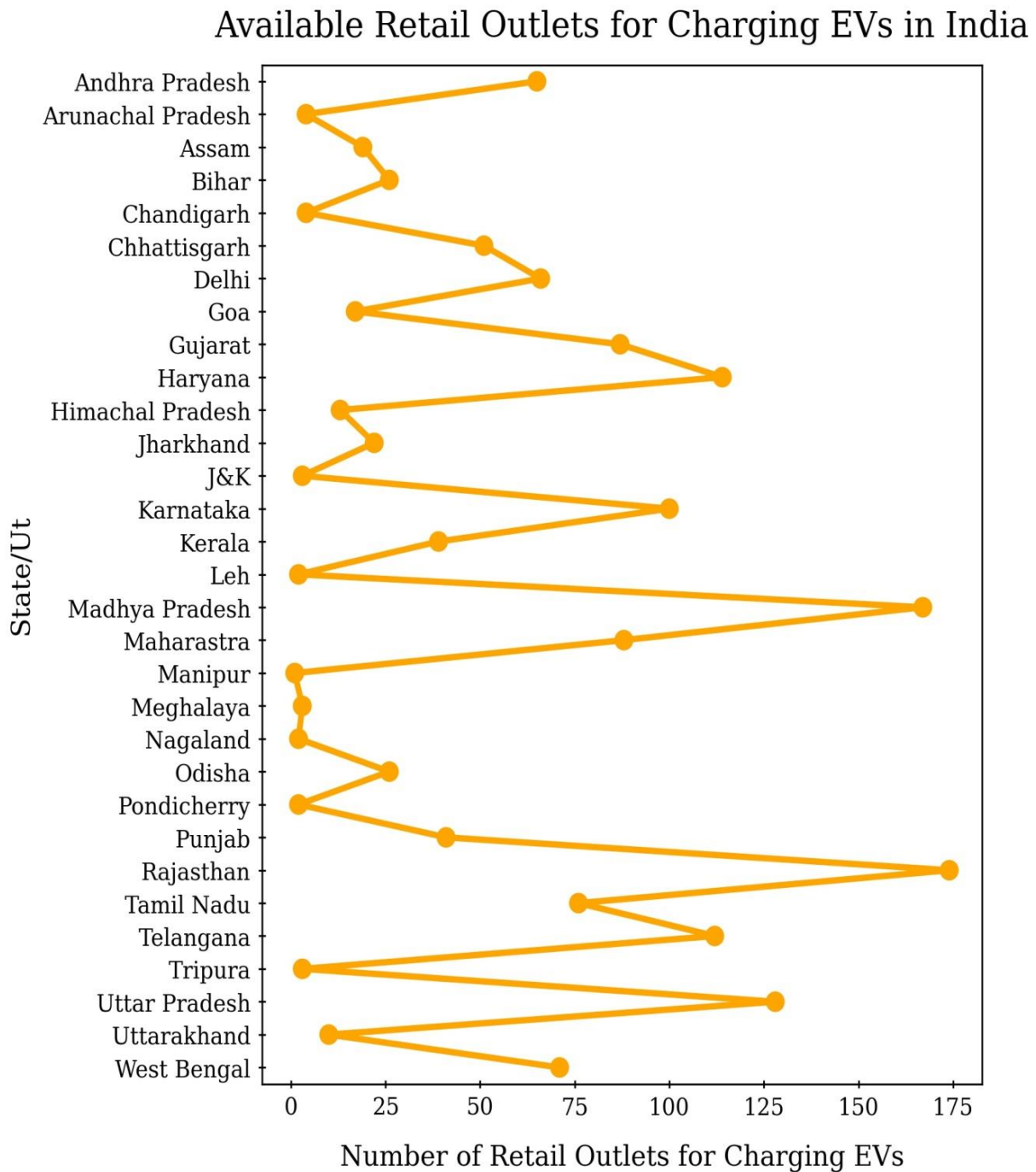
Number of 3-wheeler EVs in India



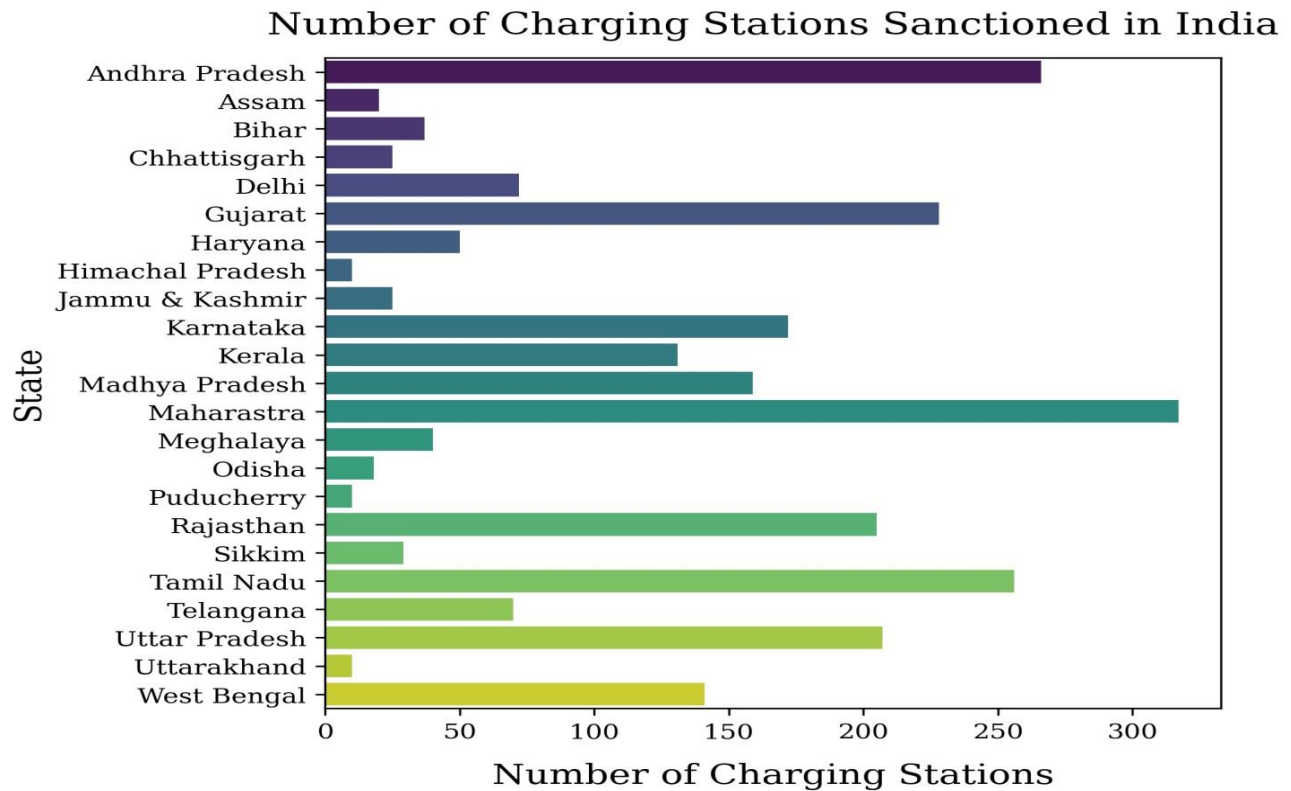
Number of 4-wheeler EVs in India



Retail outlets in India for charging EVs

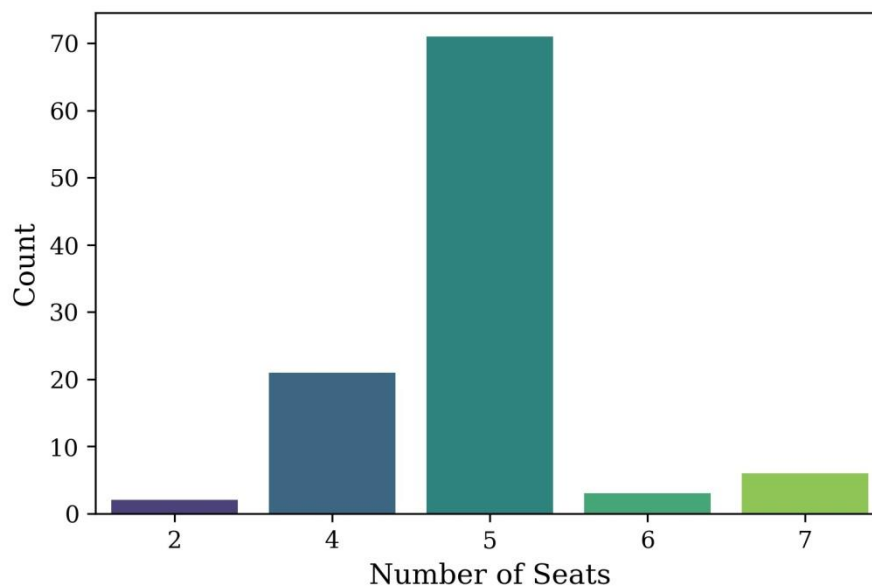


Number of charging stations sanctioned by Government of India

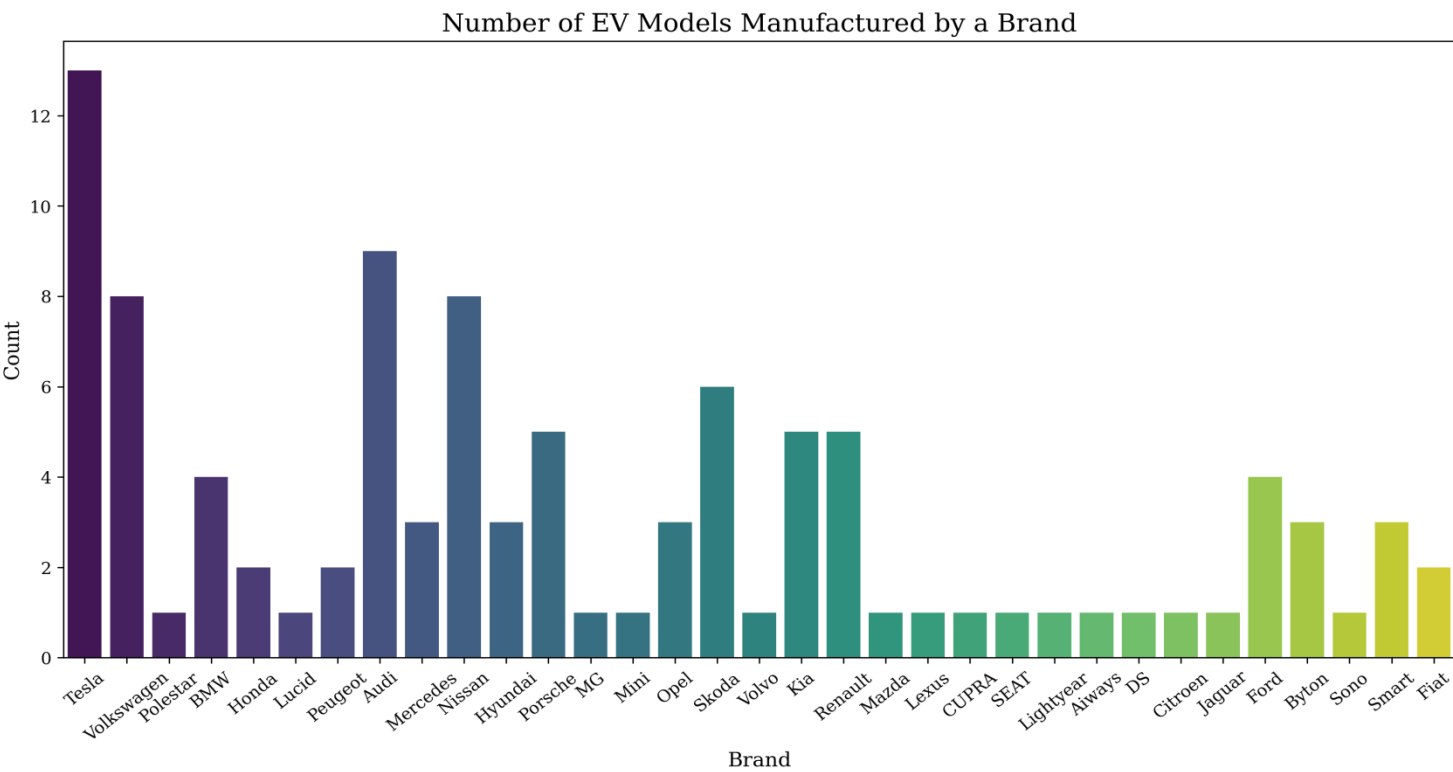


Choices for the number of seats for EVs in India

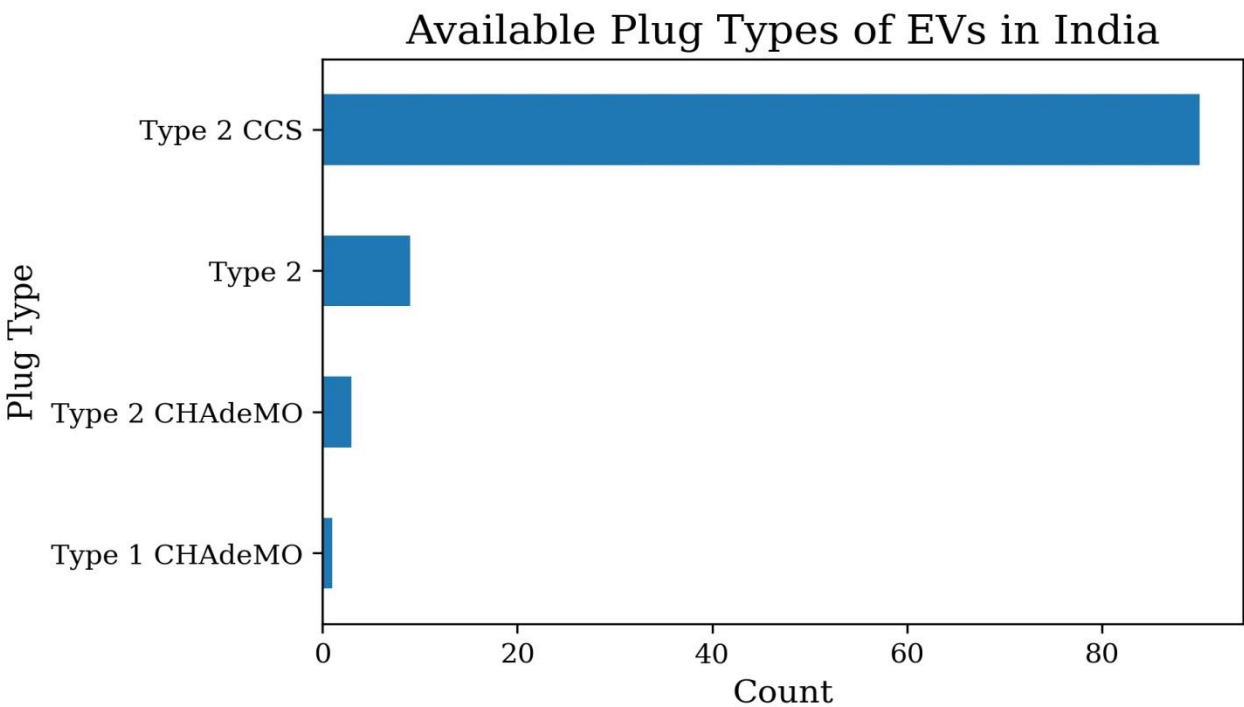
Available Electric Vehicles of Different Number of Seats in India



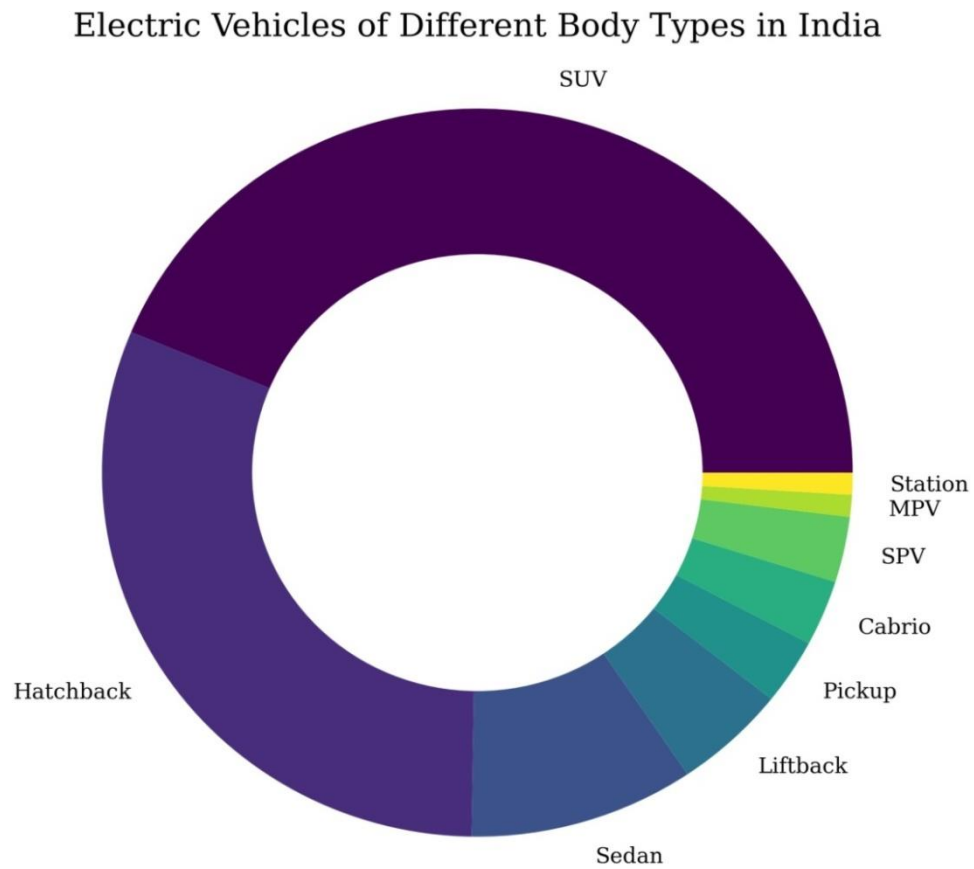
Top EV manufacturing brands in India



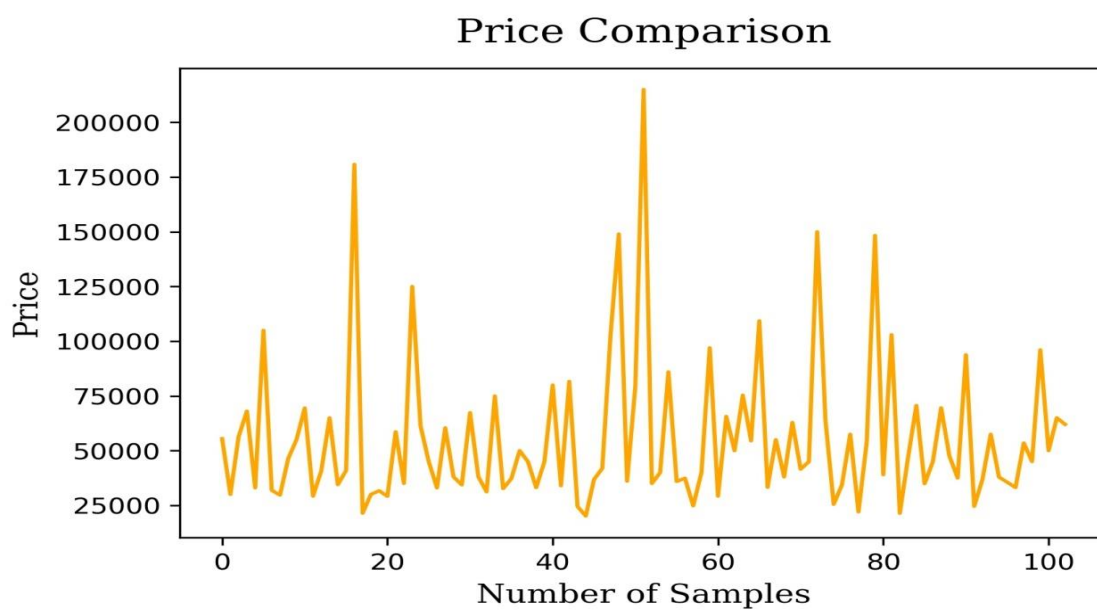
Types of EV plugs available in India



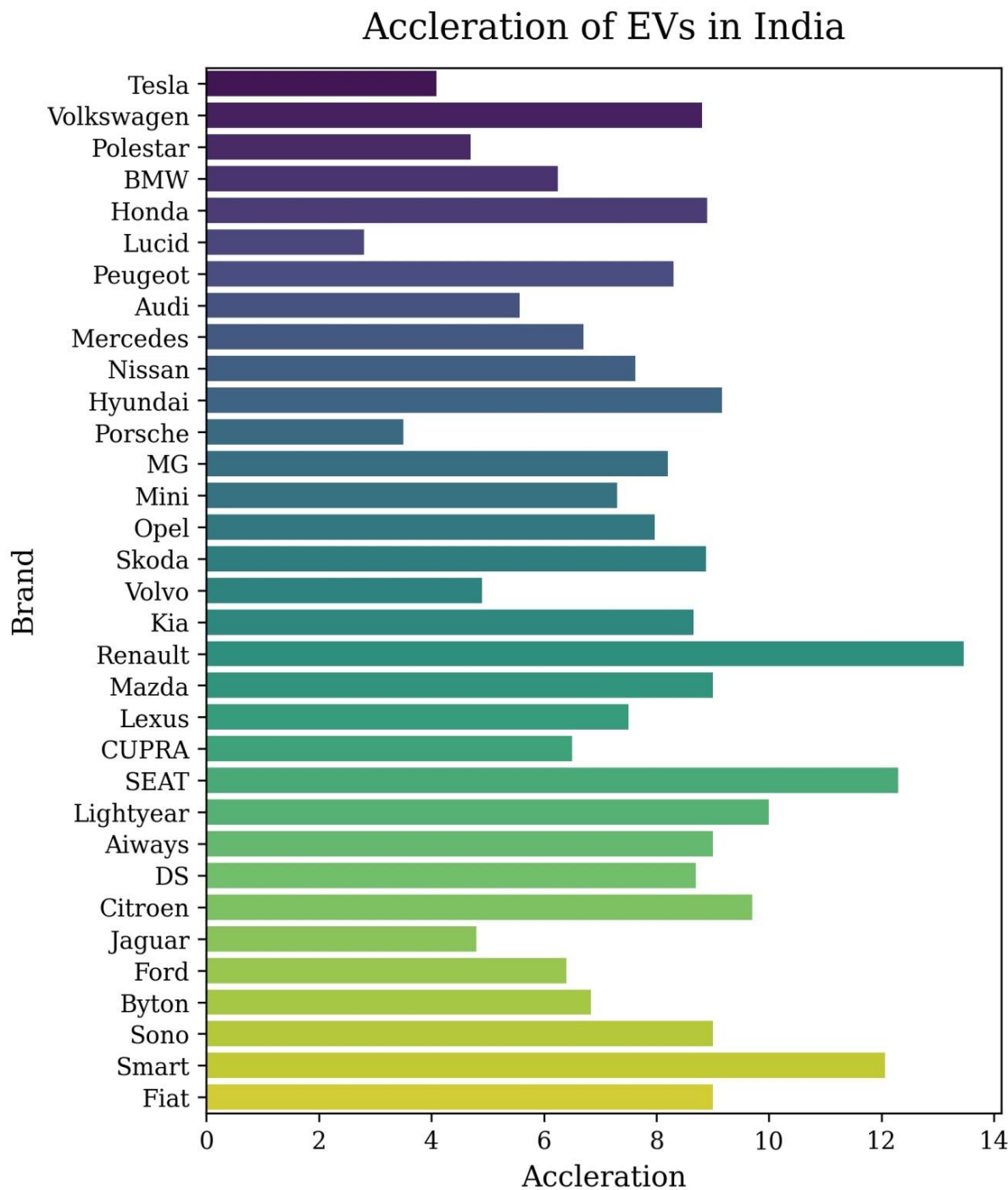
Body types of EVs in India



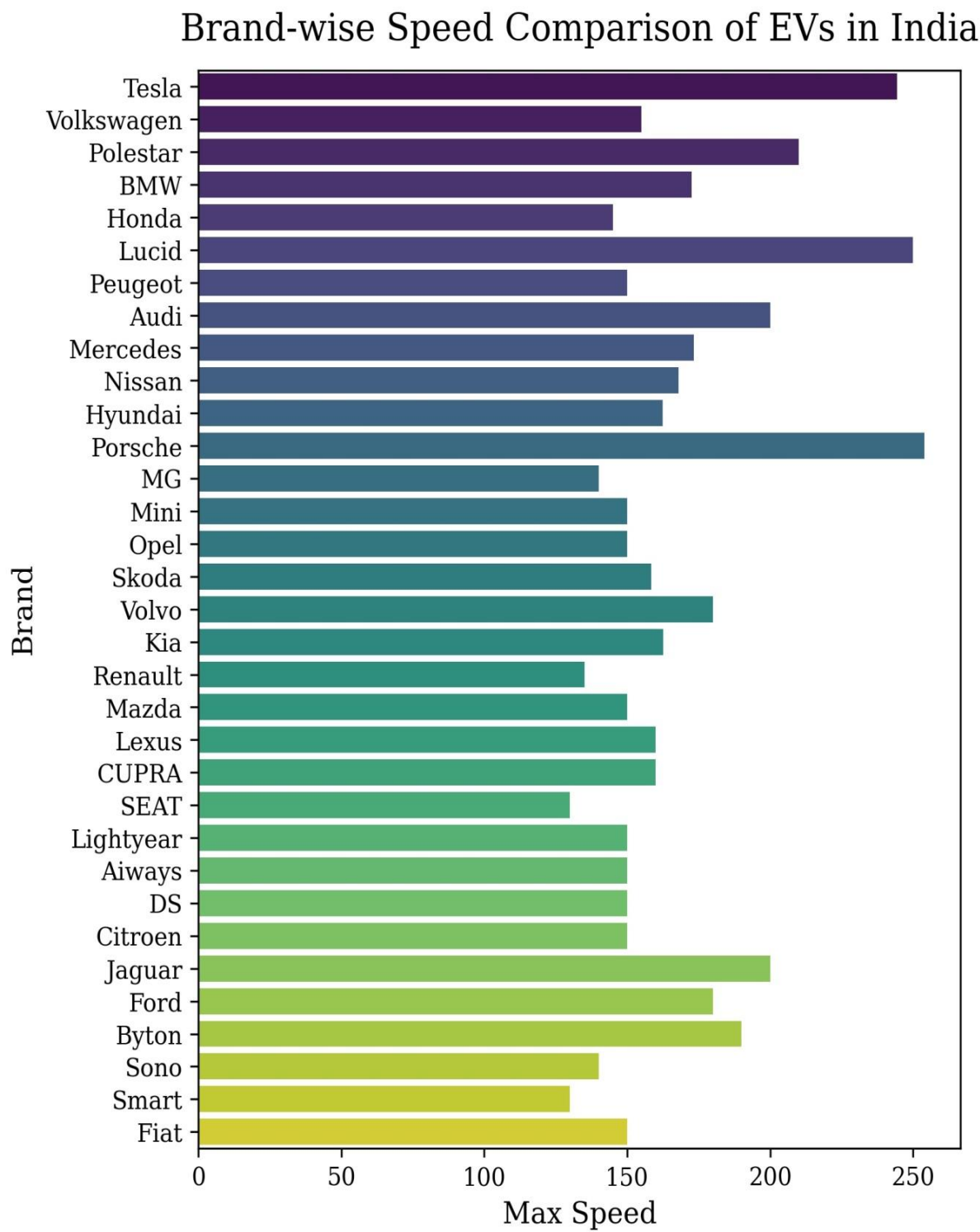
Price comparison of different brands of EVs in India



Comparison of different brands of EVs based on acceleration



Comparison of different brands of EVs based on speed



EV Segments in India

Electric Vehicles of Different Segments in India



Correlation Matrix

Correlation Matrix							
AccelSec	1.00	-0.79	-0.68	-0.38	-0.73	-0.18	-0.63
TopSpeed_KmH	-0.79	1.00	0.75	0.36	0.79	0.13	0.83
Range_Km	-0.68	0.75	1.00	0.31	0.72	0.30	0.67
Efficiency_WhKm	-0.38	0.36	0.31	1.00	0.32	0.30	0.40
FastCharge_KmH	-0.73	0.79	0.72	0.32	1.00	0.19	0.67
Seats	-0.18	0.13	0.30	0.30	0.19	1.00	0.02
PriceEuro	-0.63	0.83	0.67	0.40	0.67	0.02	1.00
AccelSec		TopSpeed_KmH	Range_Km	Efficiency_WhKm	FastCharge_KmH	Seats	PriceEuro

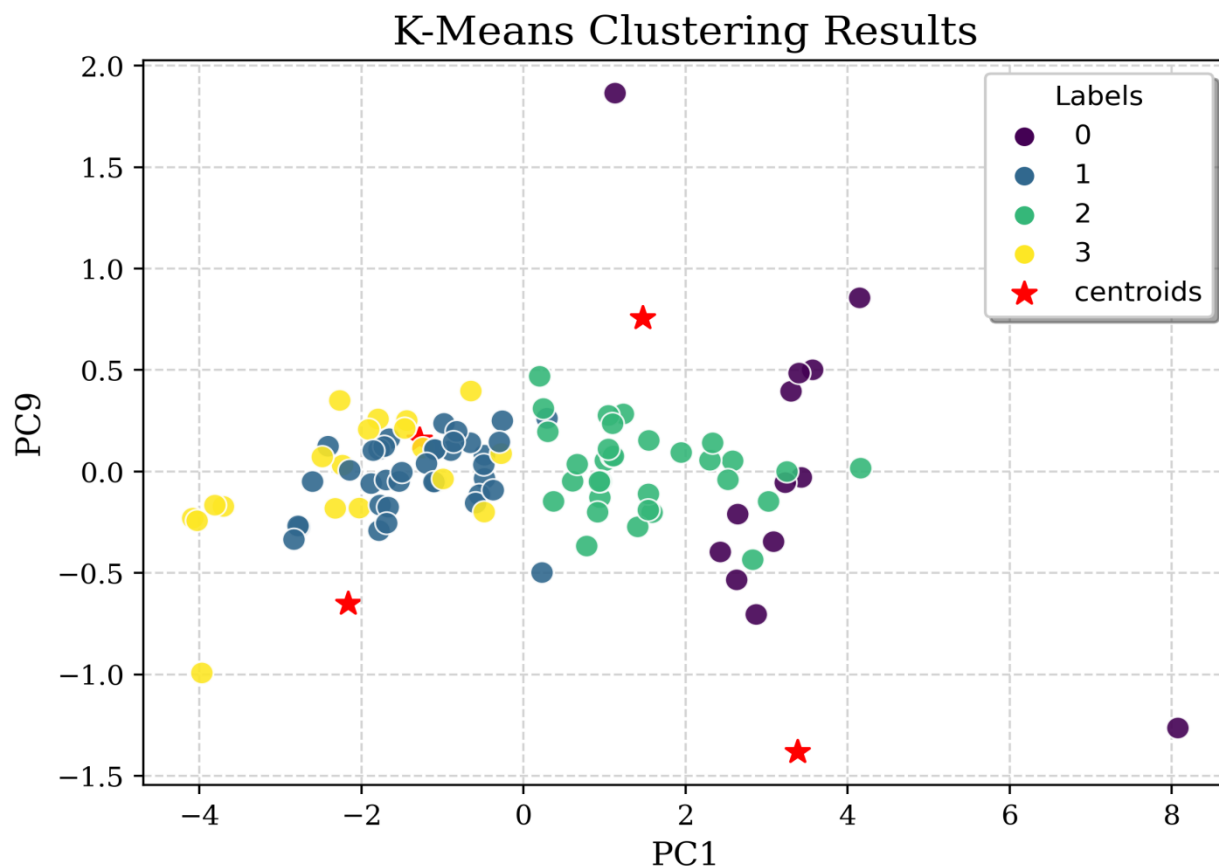
Segmentation Approaches

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).

