

# Vehicle users segmentation

This dataset has been obtained from a famous car dealer website. This dataset contains the details of around 200 Indian cars and has features like the price of the car, the type of fuel it consumes, its torque etc. Our goal here is to make user segments according to the vehicle they use.

In [31]: 

```
!importing the required packages
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sb
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.cluster import DBSCAN
from sklearn.preprocessing import OneHotEncoder
```

In [18]: 

```
# Load the datasets
id = pd.read_csv('C:\Users\ajayjuydutta\Desktop\Data analysis\Internships\Feyn Labs\Electric vehicle Ind
id.head()
```

Out[18]:

	car_name	fuel_type	engine_displacement	no_cylinder	seating_capacity	transmission_type	fuel_tank_capacity	body_type	rating	avg_price
0	Maruti Alto K10	Petrol	998	3	5	Automatic	27.0	Hatchback	4.5	49100
1	Maruti Brezza	Petrol	1462	4	5	Automatic	48.0	SUV	4.5	1097500
2	Mahindra Thar	Diesel	2184	4	4	Automatic	57.0	SUV	4.5	1478000
3	Mahindra XUV700	Diesel	2198	4	7	Automatic	60.0	SUV	4.5	1888000
4	Mahindra Scorpio-N	Diesel	2198	4	7	Automatic	57.0	SUV	4.5	1794500

In [19]: 

```
#Examining the dataset
id.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 203 entries, 0 to 202
Data columns (total 11 columns):
 #   Column              Non-Null Count  Dtype
---  ---
 0   car_name            203 non-null    object
 1   fuel_type           203 non-null    object
 2   engine_displacement 203 non-null    int64
 3   no_cylinder         203 non-null    int64
 4   seating_capacity    203 non-null    int64
 5   transmission_type   203 non-null    object
 6   fuel_tank_capacity  203 non-null    float64
 7   body_type          203 non-null    object
 8   rating              203 non-null    float64
 9   avg_price           203 non-null    float64
10   max_torque_nm       203 non-null    float64
11   max_torque_rpm      203 non-null    int64
12   max_torque_rpm_bhp 203 non-null    float64
13   max_power_hp        203 non-null    int64
dtypes: float64(4), int64(6), object(4)
memory usage: 22.9f KB
```

In [20]: 

```
#dropping null values
```

In [21]: 

```
# Correlation between ratings and average price
plt.figure(figsize=(10, 6))
plt.scatter(id['rating'], id['avg_price'], alpha=0.1)
plt.xlabel('Ratings')
plt.ylabel('Price in Rupees')
plt.title('Correlation Scatter Plot: Ratings vs. Price in rupees')
plt.grid(True)
plt.tight_layout()
plt.show()
```

Correlation Scatter Plot: Ratings vs. Price in rupees

# Calculating the correlation coefficient

```
correlation_coefficient = id['rating'].corr(id['avg_price'])
print(f"Correlation Coefficient between AQI and % Renewable: {correlation_coefficient}")
```

Correlation Coefficient: between AQI and % Renewable: 0.03350220450871702

Comment: Absolutely no correlation between rating and car price

In [22]: 

```
# correlation matrix
correlation_matrix = id.corr()
```

Out[22]:

	engine_displacement	no_cylinder	seating_capacity	fuel_tank_capacity	rating	avg_price	max_torque_nm	max_torque_rpm
engine_displacement	1.000000	0.948265	-0.277966	0.403951	-0.040786	0.632728	0.678812	0.39f
no_cylinder	0.948265	1.000000	-0.289816	0.381672	-0.089962	0.590956	0.615500	0.47f
seating_capacity	-0.277966	-0.289816	1.000000	-0.013077	-0.201402	-0.296028	-0.226848	-0.18f
fuel_tank_capacity	0.403951	0.381672	-0.013077	1.000000	0.023962	0.163430	0.094807	0.322
rating	-0.040786	-0.089962	0.201402	0.023962	1.000000	0.033502	-0.011984	-0.04f
avg_price	0.632728	0.590956	-0.296028	0.163430	0.033502	1.000000	0.636429	0.010
max_torque_nm	0.678812	0.615500	-0.226848	0.094807	-0.011984	0.636429	1.000000	0.067
max_torque_rpm	0.391929	0.474884	-0.187166	0.322906	-0.041517	0.107133	0.067723	1.000
max_power_hp	0.723541	0.686959	-0.389978	0.151109	-0.036662	0.675211	0.917634	0.84f
max_power_bhp	0.427978	0.523751	-0.209913	0.303076	-0.019821	0.090975	0.003525	0.282

In [24]: 

```
# correlation plot
plt.figure(figsize=(8, 6))
sb.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt='.2f')
```

Out[24]:

	engine_displacement	no_cylinder	seating_capacity	fuel_tank_capacity	rating	avg_price	max_torque_nm	max_torque_rpm
engine_displacement	1.00	0.95	-0.28	0.40	-0.04	0.63	0.68	0.40
no_cylinder	0.95	1.00	-0.29	0.38	-0.09	0.59	0.62	0.47
seating_capacity	-0.28	-0.29	1.00	-0.01	0.20	-0.30	-0.23	-0.19
fuel_tank_capacity	0.40	0.38	-0.01	1.00	0.02	0.16	0.09	0.32
rating	-0.04	-0.09	0.20	0.02	1.00	0.03	-0.01	-0.04
avg_price	0.63	0.59	-0.30	0.16	0.03	1.00	0.64	0.01
max_torque_nm	0.68	0.62	-0.23	0.09	-0.01	0.64	1.00	0.00
max_torque_rpm	0.40	0.47	-0.19	0.32	-0.04	0.01	0.00	1.00
max_power_hp	0.72	0.69	-0.39	0.15	-0.04	0.68	0.92	0.87
max_power_bhp	0.43	0.52	-0.21	0.33	-0.02	0.09	0.00	0.87

In [25]: 

```
# Finding a subset of variables
variables = ['engine_displacement', 'no_cylinder', 'seating_capacity', 'fuel_tank_capacity', 'rating', 'avg_price', 'max_torque_nm', 'max_torque_rpm', 'max_power_hp', 'max_power_bhp']
subset_df = id[variables]
```

Out[25]:

	engine_displacement	no_cylinder	seating_capacity	fuel_tank_capacity	rating	avg_price	max_torque_nm	max_torque_rpm	max_power_hp
0	998	3	5	27.0	4.5	49100	89.0	3500	65f
1	1462	4	5	48.0	4.5	1097500	136.8	4400	101f
2	2184	4	4	57.0	4.5	1478000	30.0	2800	190f
3	2198	4	7	60.0	4.5	1888000	45.0	2800	182f
4	2198	4	7	57.0	4.5	1794500	40.0	2750	172f
198	1991	4	5	0.0	4.5	829000	5000	5250	415f
199	1998	4	5	59.0	4.5	1041000	400.0	4400	254f
200	1956	4	7	60.0	4.5	1845000	350.0	2500	167f
201	3996	8	5	85.0	3.5	2170000	800.0	4500	591f
202	796	3	5	35.0	4.5	394000	69.0	3500	47f

In [27]: 

```
# Top 10 costliest cars
car_price = id.groupby('car_name')['avg_price'].sum()
top_ten_costliest_cars = car_price.sort_values(ascending=False).head(10)
```

Out[27]:

	car_name	avg_price
0	Lamborghini Aventador	7625000
1	Rolls-Royce Ghost	7500000
2	Jaguar i-Pace	7350000
3	Rolls-Royce Cullinan	6950000
4	Rolls-Royce Phantom	6750000
5	Ferrari 812	5750000
6	McLaren GT	4500000
7	BMW X5 M	4100000
8	Bentley Bentayga	4100000
9	Maserati GranTurismo	4020000

In [29]: 

```
# Creating a bar graph
plt.figure(figsize=(10, 6))
top_ten_costliest_cars.plot(kind='bar', color='lime')
plt.xlabel('Car model')
plt.ylabel('Average Price in Rupees')
plt.title('Top 10 cars by Price')
plt.xticks(rotation=45)
plt.grid(True)
plt.show()
```

Top 10 cars by Price

In [33]: 

```
# Creating an instance of OneHotEncoder
encoder = OneHotEncoder(sparse=False)
```

# Fitting and transforming the categorical variables

```
encoded_categories = encoder.fit_transform(id[['car_name', 'fuel_type', 'transmission_type', 'body_type']])
```

# Creating a DataFrame with the one-hot encoded columns

```
encoded_id = pd.DataFrame(encoded_categories, columns=encoder.get_feature_names(['car_name', 'fuel_type', 'transmission_type', 'body_type']))
```

Out[33]:

	car_name	fuel_type	transmission_type	body_type	rating	avg_price	max_torque_nm	max_torque_rpm	max_power_hp
0	Maruti Alto K10	Petrol	Automatic	Hatchback	4.5	49100	89.0	3500	65f
1	Maruti Brezza	Petrol	Automatic	SUV	4.5	1097500	136.8	4400	101f
2	Mahindra Thar	Diesel	Automatic	SUV	4.5	1478000	30.0	2800	190f
3	Mahindra XUV700	Diesel	Automatic	SUV	4.5	1888000	45.0	2800	182f
4	Mahindra Scorpio-N	Diesel	Automatic	SUV	4.5	1794500	40.0	2750	172f
198	Mercedes-Benz A45 S	Petrol	Automatic	Hatchback	4.5	829000	5000	5250	415f
199	BMW 3 Series Gran Limousine	Petrol	Automatic	Sedan	4.5	1041000	400.0	4400	254f
200	MG Hector Plus	Diesel	Manual	SUV	4.5	1845000	350.0	2500	167f
201	Audi RS Q8	Petrol	Automatic	SUV	3.5	2170000	800.0	4500	591f
202	Alto 800	Petrol	Manual	Hatchback	4.5	394000	69.0	3500	47f

In [35]: 

```
# Scaling the numerical features
scaler = StandardScaler()
numerical_cols = ['engine_displacement', 'no_cylinder', 'seating_capacity', 'fuel_tank_capacity', 'rating', 'avg_price', 'max_torque_nm', 'max_torque_rpm', 'max_power_hp', 'max_power_bhp']
encoded_id[numerical_cols] = scaler.fit_transform(encoded_id[numerical_cols])
```

Out[35]:

	engine_displacement	no_cylinder	seating_capacity	fuel_tank_capacity	rating	avg_price	max_torque_nm	max_torque_rpm	max_power_hp
0	-0.877790	-0.674995	-0.012792	0.663958	0.237980	-0.679006	-1.326691	0.096474	-1.04
1	-0.566332	-0.280113	-0.012792	0.064376	0.237980	-0.636945	-1.326691	0.062146	-0.85
2	-0.081818	-0.280113	-0.783335	0.376519	0.237980	-0.581057	-0.434634	-0.343494	-0.73
3	-0.072423	-0.280113	1.718335	0.048057	0.237980	-0.582123	0.199530	-0.343494	-0.40
4	-0.072423	-0.280113	1.718335	0.376519	0.237980	-0.588607	-0.011858	-0.374920	-0.48
198	0.211335	-0.280113	-0.012792	-1.600387	0.237980	-0.655665	0.410918	1.196393	0.07
199	-0.226827	-0.280113	-0.012792	0.404585	0.237980	-0.640863	-0.011858	0.621461	-0.76
200	-0.048822	-0.280113	1.718335	0.085067	0.237980	-0.585105	-0.223246	-0.523251	-0.51
201	1.135506	1.299414	-0.012792	1.347631	-0.340536	0.971862	1.679246	0.724999	1.68
202	-1.013265	-0.674995	-0.012792	-0.386497	0.237980	-0.685733	-1.411247	0.096474	-1.13

Out[37]:

	engine_displacement	no_cylinder	seating_capacity	fuel_tank_capacity	rating	avg_price	max_torque_nm	max_torque_rpm	max_power_hp
0	-0.877790	-0.674995	-0.012792	0.663958	0.237980	-0.679006	-1.326691	0.096474	-1.04
1	-0.566332	-0.280113	-0.012792	0.064376	0.237980	-0.636945	-1.326691	0.062146	-0.85
2	-0.081818	-0.280113	-0.783335	0.376519	0.237980	-0.581057	-0.434634	-0.343494	-0.73
3	-0.072423	-0.280113	1.718335	0.048057	0.237980	-0.582123	0.199530	-0.343494	-0.40
4	-0.072423	-0.280113	1.718335	0.376519	0.237980	-0.588607	-0.011858	-0.374920	-0.48
198	0.211335	-0.280113	-0.012792	-1.600387	0.237980	-0.655665	0.410918	1.196393	0.07
199	-0.226827	-0.280113	-0.012792	0.404585	0.237980	-0.640863	-0.011858	0.621461	-0.76
200	-0.048822	-0.280113	1.718335	0.085067	0.237980	-0.585105	-0.223246	-0.523251	-0.51
201	1.135506	1.299414	-0.012792	1.347631	-0.340536	0.971862	1.679246	0.724999	1.68
202	-1.013265	-0.674995	-0.012792	-0.386497	0.237980	-0.685733	-1.411247	0.096474	-1.13

In [37]: 

```
# Applying PCA to the combined numerical features
pca = PCA(n_components=2)
pca_result = pca.fit_transform(encoded_id[numerical_cols])
```

Out[37]:

	PC1	PC2
0	-1.8624445e+00	-1.0069177e+00
1	-1.0204898e+00	-1.5459305e+00
2	-7.8484952e+01	-1.0197208e+01
3	-9.7032280e+01	-1.6672961e+01
4	-1.0975732e+01	-1.7060572e+01
5	-1.6947777e+01	-1.5051802e+01
6	-1.001135e+00	-2.2914248e+01
7	-1.5812048e+00	-1.1435399e+01
8	-1.0452800e+00	-1.9572639e+01
9	-1.1838672e+00	-1.3406973e+01
10	-1.3423019e+00	-1.6755036e+01
11	-1.3257306e+00	-1.2807966e+01
12	-1.2995248e+00	-2.2865945e+01
13	-6.8356769e+01	-1.3454705e+01
14	-9.5725926e+01	-1.7619576e+01
15	-1.2426648e+00	-1.5003581e+00
16	-1.3675983e+00	9.3309746e+01
17	-1.3930288e+00	-1.5808302e+01
18	-1.2816803e+00	-2.2536245e+01
19	-1.2435285e+00	-1.5012045e+01
20	-9.1426302e+01	-1.6713743e+01
21	-2.1146028e+00	8.1655516e+01
22	-1.4687720e+00	-1.4580584e+01
23	-1.2606605e+00	-1.2831102e+01
24	-1.3502509e+00	8.3867428e+02
25	-1.4878482e+00	-1.5428484e+01
26	-1.4888108e+00	-3.1154649e+01
27	-1.2946781e+00	-2.3048308e+01
28	-1.8628484e+00	-2.1967072e+01
29	-1.9026108e+00	-1.5924846e+01
30	-1.4964057e+00	-1.1442316e+01
31	-1.7046106e+00	-1.2650770e+01
32	-9.1594863e+01	5.9830744e+01
33	-1.7133690e+00	-1.5645183e+01
34	-1.4048856e+00	2.7117939e+00
35	-1.7560230e+00	-1.6391610e+01
36	-1.4965842e+00	-1.1443789e+01
37	-1.8357476e+00	-1.0050509e+01
38	-1.3675983e+00	-2.5026175e+01
39	3.5937477e+00	6.0482164e+01
40	-1.8273567e+00	-1.6991406e+01
41	-1.9806108e+00	-1.5924846e+01
42	-1.5823654e+00	-3.2513382e+01
43	-1.4269410e+00	-1.2337022e+01
44	-8.9324842e+01	-9.5626752e+01
45	-8.9301824e+01	-9.5601766e+01
46	-1.3252032e+00	-1.3950905e+01
47	-1.3636594e+00	7.5231610e+01
48	-6.5572898e+00	2.2045610e+01
49	-8.5809278e+01	-1.0045898e+01
50	-6.3026697e+01	3.7747355e+01
51	-1.7340868e+00	-1.4593345e+01
52	-1.1614015e+00	-3.1973555e+01
53	-1.0153438e+00	-1.5414958e+01



