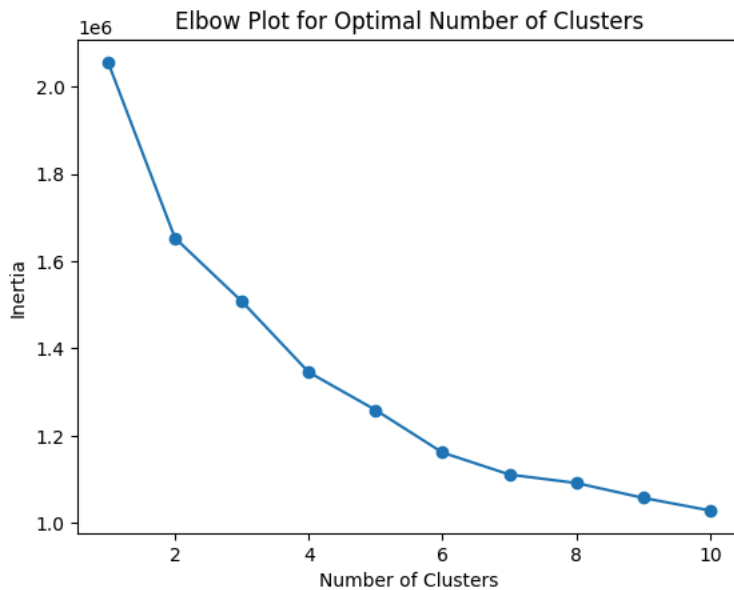


```

1 # Load necessary libraries and dataset
2 import pandas as pd
3 import numpy as np
4 from sklearn.cluster import KMeans
5 from sklearn.preprocessing import StandardScaler, OneHotEncoder
6 import matplotlib.pyplot as plt
7
8 # Load the dataset
9 file_path = "Electric_Vehicle_Population_Data 6.xlsx"
10 data = pd.read_excel(file_path)
11
12 # Define relevant columns
13 numerical_columns = ['Electric Range', 'Model Year', 'Income', 'Car Price Today',
14                     'Car Price Launched']
15 categorical_columns = ['County', 'City', 'Make', 'Model', 'Electric Vehicle Type',
16                     'Clean Alternative Fuel Vehicle (CAFV) Eligibility',
17                     'Electric Utility']
18
19 print(data.shape)
20
21 data_relevant = data[numerical_columns + categorical_columns].dropna()
22
23 print(data_relevant.shape)
24
25 # Convert all categorical columns to string type before encoding
26 # This ensures all values within a column are of the same type (string)
27 for col in categorical_columns:
28     data_relevant[col] = data_relevant[col].astype(str)
29
30 # Re-encode categorical features
31 # The 'sparse' argument has been replaced with 'sparse_output' in newer
32 # versions of scikit-learn
33 encoder = OneHotEncoder(sparse_output=False, handle_unknown='ignore')
34 encoded_categorical = encoder.fit_transform(data_relevant[categorical_columns])
35
36 # Scale numerical features
37 scaler = StandardScaler()
38 scaled_numerical = scaler.fit_transform(data_relevant[numerical_columns])
39
40 # Combine numerical and encoded categorical features
41 features_combined_reduced = np.hstack((scaled_numerical, encoded_categorical))
42
43 (204556, 20)
44 (204544, 12)
45
46 scaled_numerical.shape
47
48 (204544, 5)
49
50 features_combined_reduced.shape
51
52 (204544, 746)
53
54 # elbow plot for different cluster number
55
56 # Create an empty list to store the inertia values
57 inertia_values = []
58
59 # Define a range of cluster numbers to try
60 cluster_range = range(1, 11) # Try cluster numbers from 1 to 10
61
62 # Loop through different cluster numbers and calculate inertia
63 for n_clusters in cluster_range:
64     kmeans = KMeans(n_clusters=n_clusters, random_state=42)
65     kmeans.fit(features_combined_reduced)
66     inertia_values.append(kmeans.inertia_)
67
68 # Plot the elbow plot
69 plt.plot(cluster_range, inertia_values, marker='o')
70 plt.xlabel('Number of Clusters')
71 plt.ylabel('Inertia')
72 plt.title('Elbow Plot for Optimal Number of Clusters')
73 plt.show()

```



```

1 # Perform K-Means clustering
2 kmeans_final = KMeans(n_clusters=4, random_state=42)
3 kmeans_final.fit(features_combined_reduced)
4
5 # Assign cluster labels
6 data_relevant['Cluster'] = kmeans_final.labels_
7
8 # Analyze cluster sizes
9 cluster_counts_final = data_relevant['Cluster'].value_counts()
10
11 cluster_counts_final

```



Cluster	count
1	119005
0	33410
2	32519
3	19610

dtype: int64

```
1 kmeans_final.cluster_centers_
```



```

array([[ 2.11113334, -0.87444015,  0.07046614, ...,  0.19754565,
         0.37461838,  0.02274768],
       [-0.5300487 ,  0.57137136, -0.00818175, ...,  0.20911726,
         0.3810176 ,  0.02008319],
       [ 0.08876409, -1.59940414, -0.17172339, ...,  0.21762662,
         0.25886405,  0.03742428],
       [-0.52733499,  0.67465681,  0.2143635 , ...,  0.17526772,
         0.45048445,  0.01713412]])

```

```
1 kmeans_final.cluster_centers_.shape
```



```
(4, 746)
```

```

1 # tsne plot for these cluster centres
2
3 from sklearn.manifold import TSNE
4 import matplotlib.pyplot as plt
5
6 # Assuming kmeans_final.cluster_centers_ contains the cluster centers
7
8 # Apply t-SNE to reduce the dimensionality of the cluster centers
9 # Set perplexity to a value less than the number of cluster centers (4 in this case)
10 tsne = TSNE(n_components=2, random_state=42, perplexity=3) # perplexity should be < 4

```

```

11 cluster_centers_tsne = tsne.fit_transform(kmeans_final.cluster_centers_)
12
13 # Create a scatter plot of the reduced cluster centers
14 plt.figure(figsize=(8, 6))
15 plt.scatter(cluster_centers_tsne[:, 0], cluster_centers_tsne[:, 1], c=range(kmeans_final.n_clusters), cmap='viridis')
16 plt.xlabel('t-SNE Dimension 1')
17 plt.ylabel('t-SNE Dimension 2')
18 plt.title('t-SNE Visualization of Cluster Centers')
19 plt.colorbar()
20 plt.show()

```



```

1 # interpret all 4 cluster centres in terms of mean statistic value of each column
2
3 # Assuming 'kmeans_final.cluster_centers_' and 'numerical_columns' are defined as in your code
4
5 # Create a DataFrame from the cluster centers
6 df_cluster_centers = pd.DataFrame(kmeans_final.cluster_centers_, columns=[f"Feature_{i}" for i in range(kmeans_final.cluster_centers_.shape[1])])
7
8
9 # Interpret the cluster centers in terms of the mean statistical value of each numerical column
10
11 # Get the indices of the numerical features in the combined feature matrix
12 numerical_feature_indices = list(range(len(numerical_columns)))
13
14 # Iterate through the cluster centers and print the mean values for numerical features
15 for cluster_num in range(4):
16     print(f"\nCluster {cluster_num} Mean Values for Numerical Features:")
17     for i in numerical_feature_indices:
18         print(f" - {numerical_columns[i]}: {df_cluster_centers.iloc[cluster_num, i]}")
19
20

```

```

Cluster 0 Mean Values for Numerical Features:
 - Electric Range: 2.1111333449926826
 - Model Year: -0.8744401545444802
 - Income: 0.0704661398820383
 - Car Price Today: -0.4117284535441436
 - Car Price Launched: 0.318294442094286

```

Cluster 1 Mean Values for Numerical Features:

- Electric Range: -0.5300486997809705
- Model Year: 0.5713713610219502
- Income: -0.008181749810343357
- Car Price Today: 0.10625450918321414
- Car Price Launched: -0.2030609977787674

Cluster 2 Mean Values for Numerical Features:

- Electric Range: 0.0887640923763393
- Model Year: -1.5994041411634474
- Income: -0.17172338994117586
- Car Price Today: -1.2281677646295819
- Car Price Launched: -0.8569427024597143

Cluster 3 Mean Values for Numerical Features:

- Electric Range: -0.52733498514284
- Model Year: 0.6746568083352225
- Income: 0.21436350434565268
- Car Price Today: 2.0933104184374933
- Car Price Launched: 2.111064582947699

```

1 # Assuming kmeans_final.cluster_centers_ and numerical_columns are defined as in your code
2
3 # Create a DataFrame from the cluster centers
4 df_cluster_centers = pd.DataFrame(kmeans_final.cluster_centers_, columns=[f"Feature_{i}" for i in range(kmeans_final.cluster_centers_.shape[0])])
5
6 # Get the indices of the numerical features in the combined feature matrix
7 numerical_feature_indices = list(range(len(numerical_columns)))
8
9 # Iterate through the cluster centers and print the mean values for numerical features
10 for cluster_num in range(kmeans_final.n_clusters): # Iterate through all clusters
11     print(f"\nCluster {cluster_num} Characteristics:")
12     for i in numerical_feature_indices:
13         if numerical_columns[i] == 'Electric Range':
14             if df_cluster_centers.iloc[cluster_num, i] > 0:
15                 print(f" - Typically has a higher than average electric range.")
16             else:
17                 print(f" - Typically has a lower than average electric range.")
18
19         elif numerical_columns[i] == 'Model Year':
20             if df_cluster_centers.iloc[cluster_num, i] > 0:
21                 print(f" - Predominantly consists of more recent model year vehicles.")
22             else:
23                 print(f" - Predominantly consists of older model year vehicles.")
24
25         elif numerical_columns[i] == 'Income':
26             if df_cluster_centers.iloc[cluster_num, i] > 0:
27                 print(f" - Associated with higher than average income levels.")
28             else:
29                 print(f" - Associated with lower than average income levels.")
30
31         elif numerical_columns[i] == 'Car Price Today':
32             if df_cluster_centers.iloc[cluster_num, i] > 0:
33                 print(f" - Tends to have a higher than average current car price.")
34             else:
35                 print(f" - Tends to have a lower than average current car price.")
36
37         elif numerical_columns[i] == 'Car Price Launched':
38             if df_cluster_centers.iloc[cluster_num, i] > 0:
39                 print(f" - Tends to have had a higher than average launch price.")
40             else:
41                 print(f" - Tends to have had a lower than average launch price.")

```



Cluster 0 Characteristics:

- Typically has a higher than average electric range.

- Predominantly consists of older model year vehicles.
- Associated with higher than average income levels.
- Tends to have a lower than average current car price.
- Tends to have had a higher than average launch price.

#### Cluster 1 Characteristics:

- Typically has a lower than average electric range.
- Predominantly consists of more recent model year vehicles.
- Associated with lower than average income levels.
- Tends to have a higher than average current car price.
- Tends to have had a lower than average launch price.

#### Cluster 2 Characteristics:

- Typically has a higher than average electric range.
- Predominantly consists of older model year vehicles.
- Associated with lower than average income levels.
- Tends to have a lower than average current car price.
- Tends to have had a lower than average launch price.

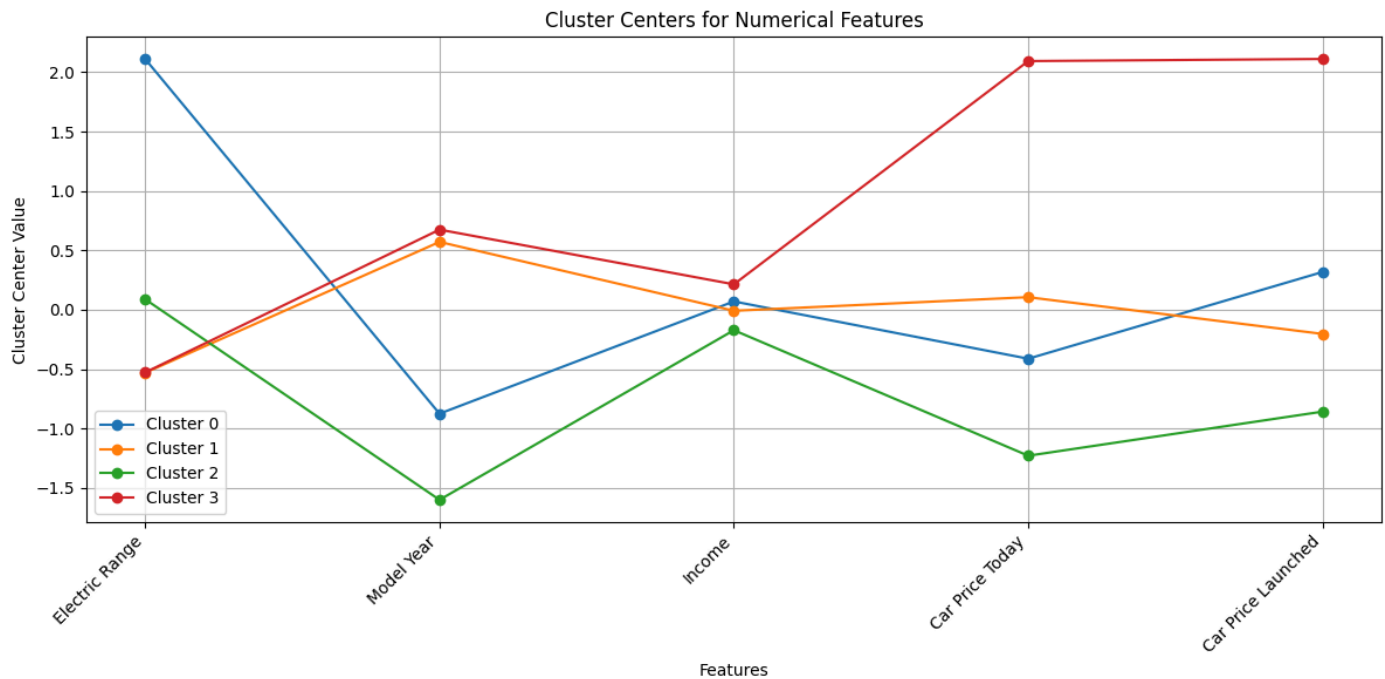
#### Cluster 3 Characteristics:

- Typically has a lower than average electric range.
- Predominantly consists of more recent model year vehicles.
- Associated with higher than average income levels.
- Tends to have a higher than average current car price.
- Tends to have had a higher than average launch price.

```

1 # Create a DataFrame from the cluster centers
2 df_cluster_centers = pd.DataFrame(kmeans_final.cluster_centers_, columns=[f"Feature_{i}" for i in range(kmeans_final.cluster_centers_.shape[0])]
3
4 # Get the indices of the numerical features in the combined feature matrix
5 numerical_feature_indices = list(range(len(numerical_columns)))
6
7 # Create a plot with cluster centers
8 plt.figure(figsize=(12, 6))
9 for cluster_num in range(4):
10     plt.plot(df_cluster_centers.iloc[cluster_num, numerical_feature_indices], marker='o', label=f'Cluster {cluster_num}')
11
12 plt.xticks(range(len(numerical_columns)), numerical_columns, rotation=45, ha="right")
13 plt.xlabel("Features")
14 plt.ylabel("Cluster Center Value")
15 plt.title("Cluster Centers for Numerical Features")
16 plt.legend()
17 plt.grid(True)
18 plt.tight_layout()
19 plt.show()

```



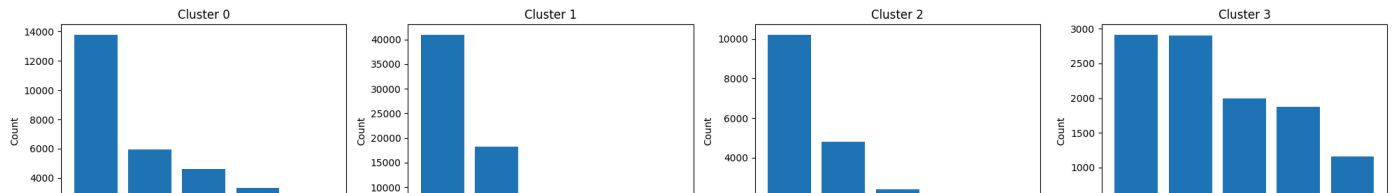
```

1 # comparative plot for non numerical columns and value count of top 5 models in each of the clusters
2
3 import matplotlib.pyplot as plt
4
5 # Assuming data_relevant and categorical_columns are defined as in your code
6
7 # Function to plot top 5 models in each cluster
8 def plot_top_models_per_cluster(data, cluster_column, categorical_column):
9     """Plots the top 5 models within each cluster."""
10
11     fig, axs = plt.subplots(1, 4, figsize=(20, 5))
12     fig.suptitle(f"Top 5 {categorical_column} in Each Cluster", fontsize=16)
13
14     for cluster_num in range(4):
15         cluster_data = data[data[cluster_column] == cluster_num]
16         top_models = cluster_data[categorical_column].value_counts().head(5)
17         axs[cluster_num].bar(top_models.index, top_models.values)
18         axs[cluster_num].set_title(f"Cluster {cluster_num}")
19         axs[cluster_num].tick_params(axis='x', rotation=45, labelsz=8)
20         axs[cluster_num].set_xlabel(f"{categorical_column}")
21         axs[cluster_num].set_ylabel("Count")
22
23     plt.tight_layout()
24     plt.show()
25
26 # Example: Plot top 5 models in each cluster
27 plot_top_models_per_cluster(data_relevant, 'Cluster', 'Model')
28

```



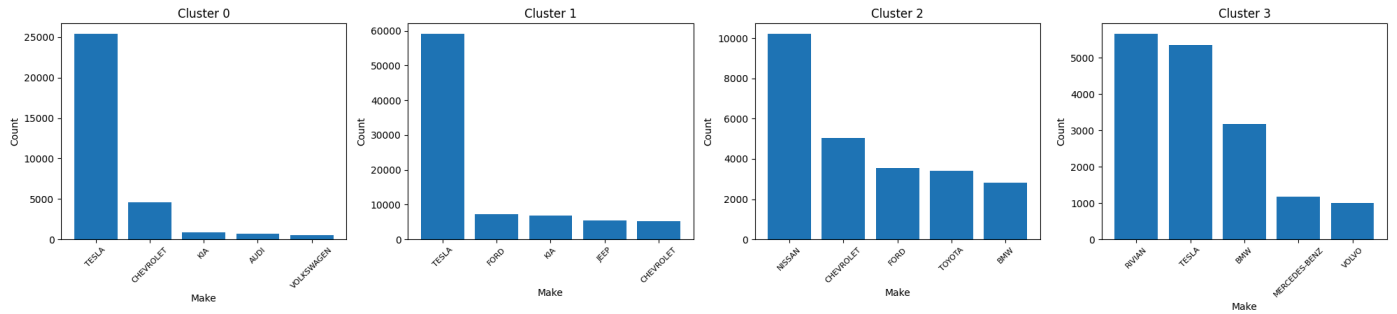
Top 5 Model in Each Cluster



```
1 # Example: Plot top 5 make in each cluster
2 plot_top_models_per_cluster(data_relevant, 'Cluster', 'Make')
```



Top 5 Make in Each Cluster



```
1 # Example: Plot top 5 city in each cluster
2 plot_top_models_per_cluster(data_relevant, 'Cluster', 'City')
```



Top 5 City in Each Cluster

