

## ▼ Problem Statement

Program to implement k-means clustering technique using any standard dataset available in the public domain

---

## ▼ Dataset Description

In this project, we will be using the dataset holding the information of carbon dioxide emission from different car models.

The dataset includes 36 instances with 5 columns which can be briefed as:

| Column | Description  |
|--------|--|
| Car    | Brand of the car   |
| Model  | Model of the car   |
| Volume | Total space available inside the car (in <i>litres</i> ) |
| Weight | Total weight of the car (in <i>kg</i> )                  |
| $CO_2$ | Total emission of carbon dioxide from the car            |

**Note:** *(This is a manually created custom dataset for this project.)*

---

## ▼ List of Activities

**Activity 1:** Import Modules and Read Data

**Activity 2:** Data Cleaning

**Activity 3:** Find Optimal Value of  $k$

**Activity 4:** Plot Silhouette Scores

## ▼ Activity 1: Import Modules and Read Data

Import the necessary Python modules along with the following modules:

- `KMeans` - For clustering using K-means.
- `re` - To remove unwanted rows using regex.

Read the data from a CSV file to create a Pandas DataFrame and go through the necessary data-cleaning process (if required).

```
1 # Import the modules and Read the data.
2 import numpy as np
3 import pandas as pd
4 df=pd.read_csv('https://raw.githubusercontent.com/jiss-sngce/CO_3/main/jkcars.csv')
5
6 # Print the first five records
7 df.head()
8
```

|   | Car        | Model      | Volume | Weight | CO2 |
|---|------------|------------|--------|--------|-----|
| 0 | Mitsubishi | Space Star | 1200   | 1160   | 95  |
| 1 | Skoda      | Citigo     | 1000   | 929    | 95  |
| 2 | Fiat       | 500        | 900    | 865    | 90  |
| 3 | Mini       | Cooper     | 1500   | 1140   | 105 |
| 4 | VW         | Up!        | 1000   | 929    | 105 |

```
1 # Get the total number of rows and columns, data types of columns and missing values (i
2 df.shape
3 df.dtypes
4 df.isnull().sum()
```

```
Car      0
Model    0
Volume   0
Weight   0
CO2      0
dtype: int64
```

### ▼ Activity 3: Find Optimal value of K

In this activity, you need to find the optimal value of  $K$  using the silhouette score.

1. Create a subset of the dataset consisting of three columns i.e Volume, Weight, and CO2.

```
1 # Create a new DataFrame consisting of three columns 'Volume', 'Weight', 'CO2'.
2 new_df = df[['Volume','Weight','CO2']]
3
4 # Print the first 5 rows of this new DataFrame.
5 new_df.head()
```

|   | Volume | Weight | CO2 |
|---|--------|--------|-----|
| 0 | 1200   | 1160   | 95  |
| 1 | 1000   | 929    | 95  |
| - | -      | -      | -   |

2. Compute K-Means clustering for the 3D dataset `data_3d` by varying  $k$  from 2 to 10 clusters. Also, for each  $k$ , calculate silhouette score using `silhouette_score` function.

### Steps to Follow

- Create an empty list to store silhouette scores obtained for each  $k$  (let's say `sil_scores`).
- Initiate a `for` loop that ranges from 2 to 10.
- Perform K-means clustering for the current value of  $k$  inside `for` loop.
- Use `fit()` and `predict()` to create clusters.
- Calculate silhouette score for current  $k$  value using `silhouette_score()` function and append it to the empty list `sil_scores`.
- Create a DataFrame with two columns. The first column must contain  $k$  values from 2 to 10 and the second column must contain silhouette values obtained after the `for` loop.

```

1 #Calculate inertia for different values of 'K'.
2 from sklearn.metrics import silhouette_score
3 from sklearn.cluster import KMeans
4 # Create an empty list to store silhouette scores obtained for each 'K'
5 sil_scores = []
6 clusters = range(2,11)
7
8 for k in clusters:
9     kmean_k = KMeans (n_clusters = k, random_state=10)
10    kmean_k.fit(new_df)
11    cluster_labels = kmean_k.predict(new_df)
12    sil_scores.append(silhouette_score(new_df,cluster_labels))
13
14 sil_data = pd.DataFrame({'K value':clusters,'silhouette_score':sil_scores})
15 sil_data
16 #sil_scores

```

|   | K value | silhouette_score |
|---|---------|------------------|
| 0 | 2       | 0.466982         |
| 1 | 3       | 0.569304         |
| 2 | 4       | 0.506027         |
| 3 | 5       | 0.537547         |

**Q:** What are the maximum silhouette score and the corresponding cluster value?

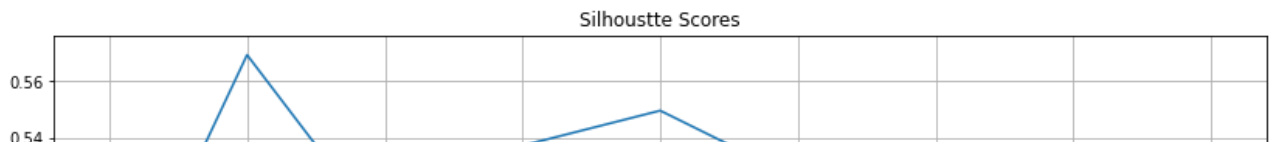
**A:**

- Maximum silhouette score: 0.569304
  - Corresponding cluster value: 3
- 

#### ▼ Activity 4: Plot silhouette Scores & WCSS Scores to find optimal value for K

Create a line plot with  $K$  ranging from 2 to 10 on the  $x$ -axis and the silhouette scores stored in `sil_scores` list on the  $y$ -axis.

```
1 # Plot silhouette scores vs number of clusters.
2 import matplotlib.pyplot as plt
3 import numpy as np
4 plt.figure(figsize=(14,5))
5 plt.title('Silhoustte Scores')
6 x = np.arange(2,11)
7 plt.xlabel('K')
8 plt.ylabel('Silhoustte scores')
9 plt.grid()
10 y = (sil_scores)
11 plt.plot(x,y)
12 plt.show()
13
```



**Q:** Write your observations of the graph.

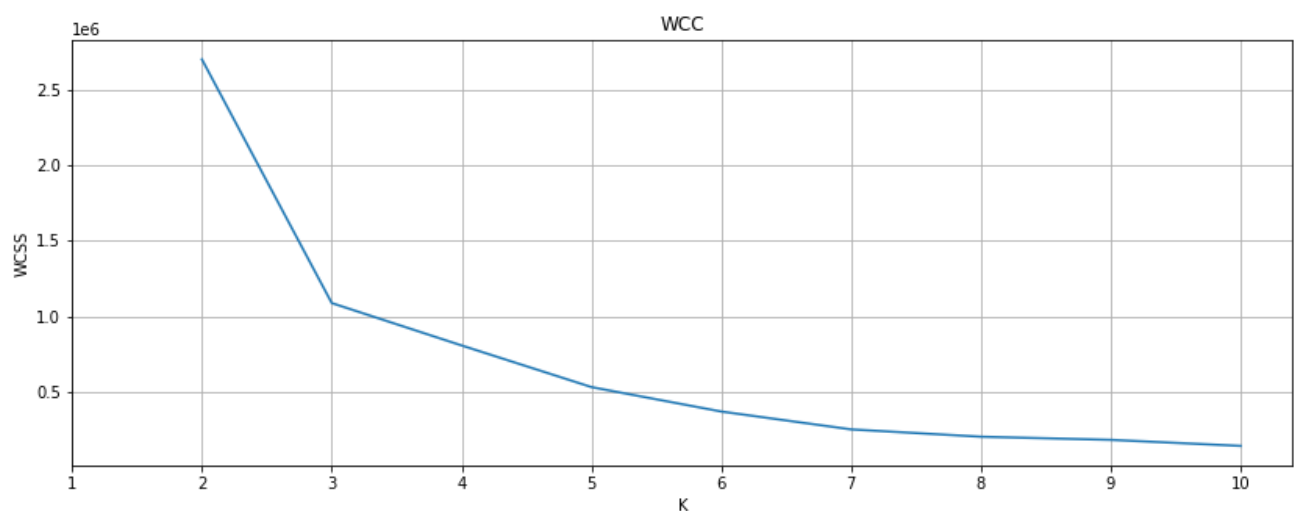
**A:** From the graph, we can conclude that the optimal value of  $k$  is 3.



```

1 # S3.1: Determine 'K' using Elbow method.
2 from sklearn.metrics import silhouette_score
3 from sklearn.cluster import KMeans
4 wcss=[]
5 clusters=range(2,11)
6
7
8 # Initiate a for loop that ranges from 1 to 10.
9 for k in clusters:
10     # Inside for loop, perform K-means clustering for current value of K. Use 'fit()' t
11     kmean_k=KMeans(n_clusters=k,random_state=10)
12     kmean_k.fit(new_df)
13
14     # Find wcss for current K value using 'inertia_' attribute and append it to the emp
15     wcss.append(kmean_k.inertia_)
16
17 # Plot WCSS vs number of clusters.
18 plt.figure(figsize=(14,5))
19 plt.title('WCC')
20 plt.plot(clusters,wcss)
21 plt.xlabel('K')
22 plt.ylabel("WCSS")
23 plt.grid()
24 plt.xticks(range(1,11))
25 plt.show()
26

```



```
1 # Clustering the dataset for K = 3
2 from sklearn.cluster import KMeans
3
4 # Perform K-Means clustering with n_clusters = 4 and random_state = 10
5 kmeans_model=KMeans(n_clusters=3,random_state=10)
6
7
8 # Fit the model to the scaled_df
9 kmeans_model.fit(new_df)
10
11
12 # Make a series using predictions by K-Means
13 cluster_labels=pd.Series(kmeans_model.predict(new_df))
14 cluster_labels.value_counts()
```

```
2    16
1     9
0     7
dtype: int64
```

```
1 # Create a DataFrame with cluster labels for cluster visualisation
2 km_df=pd.concat([df,cluster_labels],axis=1)
3 km_df.columns=list(df.columns)+['label']
4 km_df
5
```

|    | Car        | Model      | Volume | Weight | CO2 | label |
|----|------------|------------|--------|--------|-----|-------|
| 0  | Mitsubishi | Space Star | 1200   | 1160   | 95  | 0     |
| 1  | Skoda      | Citigo     | 1000   | 929    | 95  | 0     |
| 2  | Fiat       | 500        | 900    | 865    | 90  | 0     |
| 3  | Mini       | Cooper     | 1500   | 1140   | 105 | 2     |
| 4  | VW         | Up!        | 1000   | 929    | 105 | 0     |
| 5  | Skoda      | Fabia      | 1400   | 1109   | 90  | 2     |
| 6  | Ford       | Fiesta     | 1500   | 1112   | 98  | 2     |
| 7  | Audi       | A1         | 1600   | 1150   | 99  | 2     |
| 8  | Hyundai    | I20        | 1100   | 980    | 99  | 0     |
| 9  | Suzuki     | Swift      | 1300   | 990    | 101 | 0     |
| 10 | Ford       | Fiesta     | 1000   | 1112   | 99  | 0     |
| 11 | Honda      | Civic      | 1600   | 1252   | 94  | 2     |
| 12 | Hundai     | I30        | 1600   | 1326   | 97  | 2     |
| 13 | Opel       | Astra      | 1600   | 1330   | 97  | 2     |
| 14 | BMW        | 1          | 1600   | 1365   | 99  | 2     |
| 15 | Mercedes   | 2          | 2200   | 1280   | 104 | 1     |
| 16 | Volvo      | 40         | 1600   | 1415   | 109 | 2     |
| 17 | Ford       | Focus      | 2000   | 1328   | 105 | 1     |
| 18 | Ford       | Mondeo     | 1600   | 1584   | 94  | 2     |
| 19 | Mercedes   | C-Class    | 2100   | 1365   | 99  | 1     |
| 20 | Skoda      | Octavia    | 1600   | 1415   | 99  | 2     |
| 21 | Volvo      | S60        | 2000   | 1415   | 99  | 1     |
| 22 | Mercedes   | CLA        | 1500   | 1465   | 102 | 2     |
| 23 | Audi       | A4         | 2000   | 1490   | 104 | 1     |
| 24 | Audi       | A6         | 2000   | 1725   | 114 | 1     |
| 25 | Volvo      | V70        | 1600   | 1523   | 109 | 2     |
| 26 | BMW        | 5          | 2000   | 1705   | 114 | 1     |
| 27 | Volvo      | XC70       | 2000   | 1746   | 117 | 1     |
| 28 | Ford       | B-Max      | 1600   | 1235   | 104 | 2     |
| 29 | BMW        | 216        | 1600   | 1390   | 108 | 2     |